

The Telecommunications and Data Acquisition Progress Report 42-110

April-June 1992

E. C. Posner
Editor

(NASA-CR-191228) THE
TELECOMMUNICATIONS AND DATA
ACQUISITION REPORT Progress Report,
Apr. - Jun. 1992 (JPL) 290 p

N93-19413
--THRU--
N93-19436
Unclas

G3/32 0128433

August 15, 1992



National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California



The Telecommunications and Data Acquisition Progress Report 42-110

April–June 1992

E. C. Posner
Editor

August 15, 1992

NASA

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

The research described in this publication was carried out by the Jet Propulsion Laboratory, California Institute of Technology, under a contract with the National Aeronautics and Space Administration.

Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California Institute of Technology.

Preface

This quarterly publication provides archival reports on developments in programs managed by JPL's Office Telecommunications and Data Acquisition (TDA). In space communications, radio navigation, radio science, and ground-based radio and radar astronomy, it reports on activities of the Deep Space Network (DSN) in planning, in supporting research and technology, in implementation, and in operations. Also included is standards activity at JPL for space data and information systems and reimbursable DSN work performed for other space agencies through NASA. The preceding work is all performed for NASA's Office of Space Operations (OSO). The TDA Office also performs work funded by two other NASA program offices through and with the cooperation of the Office of Space Operations. These are the Orbital Debris Radar Program (with the Office of Space Station) and 21st Century Communication Studies (with the Office of Aeronautics and Exploration Technology).

In the search for extraterrestrial intelligence (SETI), the *TDA Progress Report* reports on implementation and operations for searching the microwave spectrum. In solar system radar, it reports on the uses of the Goldstone Solar System Radar for scientific exploration of the planets, their rings and satellites, asteroids, and comets. In radio astronomy, the areas of support include spectroscopy, very long baseline interferometry, and astrometry. These three programs are performed for NASA's Office of Space Science and Applications (OSSA) with the Office of Space Operations for funding DSN operational support.

Finally, tasks funded under the JPL Director's Discretionary Fund and the Caltech President's Fund which involve the TDA Office are included.

This and each succeeding issue of the *TDA Progress Report* will present material in some, but not necessarily all, for the following categories:

OSO Tasks:

DSN Advanced Systems

- Tracking and Ground-Based Navigation
- Communications, Spacecraft-Ground
- Station Control and System Technology
- Network Data Processing and Productivity

DSN Systems Implementation

- Capabilities for Existing Projects
- Capabilities for New Projects
- New Initiatives
- Network Upgrade and Sustaining

DSN Operations

- Network Operations and Operations Support
- Mission Interface and Support
- TDA Program Management and Analysis
- Ground Communications Implementation and Operations
- Data and Information Systems
- Flight-Ground Advanced Engineering
- Long-Range Program Planning

OSO Cooperative Tasks:

- Orbital Debris Radar Program
- 21st Century Communication Studies

OSSA Tasks:

Search for Extraterrestrial Intelligence

Goldstone Solar System Radar

Radio Astronomy

Discretionary Funded Tasks

Contents

OSO TASKS DSN Advanced Systems TRACKING AND GROUND-BASED NAVIGATION

Precise Tracking of the Magellan and Pioneer Venus Orbiters by Same-Beam Interferometry, Part I: Data Accuracy Analysis	1-1
J. S. Border, W. M. Folkner, R. D. Kahn, and K. S. Zukor NASA Code 310 10 60 91 01	
Application of High-Precision Two-Way Ranging to Galileo Earth-1 Encounter Navigation	21-2
V. M. Pollmeier and S. W. Thurman NASA Code 310 10 63 84 02	
The Effect of Tropospheric Fluctuations on the Accuracy of Water Vapor Radiometry	33-3
J. Z. Wilcox NASA Code 310 10 60 86 02	
The Goldstone Real-Time Connected Element Interferometer	52-4
C. Edwards, Jr., D. Rogstad, D. Fort, L. White, and B. Iijima NASA Code 310 10 63 88 01	
Spacecraft-Spacecraft Very Long Baseline Interferometry, Part I: Error Modeling and Observable Accuracy	63-5
C. Edwards, Jr., and J. S. Border NASA Code 310 10 60 91 01	
Use of the VLBI Delay Observable for Orbit Determination of Earth-Orbiting VLBI Satellites	77-6
J. S. Ulvestad NASA Code 310 10 63 50 00	
Modeling the Global Positioning System Signal Propagation Through the Ionosphere	92-7
S. Bassiri and G. Hajj NASA Code 310 10 61 84 02	
Evaluation of the Table Mountain Ronchi Telescope for Angular Tracking	104-8
G. Lanyi, G. Purcell, R. Treuhart, and A. Buffington NASA Code 310 10 60 90 03	
Deep-Space Navigation Applications of Improved Ground-Based Optical Astrometry	118-9
G. W. Null, W. M. Owen, Jr., and S. P. Synnott NASA Code 310 10 63 55 00	

COMMUNICATIONS, SPACECRAFT-GROUND

A Method for Modeling Discontinuities in a Microwave Coaxial Transmission Line	128-10
T. Y. Otoshi NASA Code 310 20 65 86 08	
DSS-13 Beam Waveguide Antenna Frequency Stability	151-11
T. Y. Otoshi and M. M. Franco NASA Code 310 20 65 92 01	
Locally Adaptive Vector Quantization: Data Compression With Feature Preservation	163-12
K.-M. Cheung and M. Sayano NASA Code 310 30 71 83 02	
Data Compression by Wavelet Transforms	179-13
M. Shahshahani NASA Code 310 30 71 83 02	

Maximal Codeword Lengths in Huffman Codes	188	-14
Y. S. Abu-Mostafa and R. J. McEliece NASA Code 310 30 71 83 02		
Comparisons of Theoretical Limits for Source Coding With Practical Compression Algorithms	194	-15
F. Pollara and S. Dolinar NASA Code 310 30 71 83 02		
Cascaded Convolutional Codes	202	-16
F. Pollara and D. Divsalar NASA Code 310 30 71 83 02		
Binary Weight Distributions of Some Reed-Solomon Codes	208	-17
F. Pollara and S. Arnold NASA Code 310 30 71 83 04		
Multiple Symbol Partially Coherent Detection of MPSK	216	-18
M. K. Simon and D. Divsalar NASA Code 310 30 71 86 99		

STATION CONTROL AND SYSTEM TECHNOLOGY

A Parameter and Configuration Study of the DSS-13 Antenna Drives	231	-19
W. Gawronski and J. A. Mellstrom NASA Code 310 40 41 10 15		

DSN Systems Implementation Capabilities for New Projects

SETI Low-Frequency Feed Design Study for DSS 24	246	-20
P. H. Stanton and P. R. Lee NASA Code 314 30 69 40 06		

NETWORK UPGRADE AND SUSTAINING

Feed-Forward Control Upgrade of the Deep Space Network Antennas	253	-21
W. Gawronski NASA Code 644 11 00 06 13		

DSN Operations NETWORK OPERATIONS AND OPERATIONS SUPPORT

Availability Analysis of the Traveling-Wave Maser Amplifiers in the Deep Space Network, Part I: The 70-Meter Antennas	263	-22
T. N. Issa NASA Code 314 40 21 30 02		

TDA PROGRAM MANAGEMENT AND ANALYSIS

A Model for the Cost of Doing a Cost Estimate	278
D. S. Remer and H. R. Buchanan NASA Code 314 40 31 30 03	
References	284

N93-19414

51-13

128434

P-20

Precise Tracking of the Magellan and Pioneer Venus Orbiters by Same-Beam Interferometry

Part I: Data Accuracy Analysis

J. S. Border, W. M. Folkner, R. D. Kahn, and K. S. Zukor
Tracking Systems and Applications Section

Simultaneous tracking of two spacecraft in orbit about a distant planet by two widely separated Earth-based radio antennas provides more-accurate positioning information than can be obtained by tracking each spacecraft separately. A demonstration of this tracking technique, referred to as same-beam interferometry (SBI), is currently being done using the Magellan and Pioneer 12 orbiters at Venus. Signals from both spacecraft fall within the same beamwidth of the Deep Space Station antennas. The plane-of-sky position difference between spacecraft is precisely determined by doubly differenced phase measurements. This radio metric measurement naturally complements line-of-sight Doppler. Data were first collected from Magellan and Pioneer 12 on August 11-12, 1990, shortly after Magellan was inserted into Venus orbit. Data were subsequently acquired in February and April 1991, providing a total of 34 hours of same-beam radio metric observables. Same-beam radio metric residuals have been analyzed and compared with model measurement error predictions. The predicted error is dominated by solar plasma fluctuations. The rms of the residuals is less than predicted by about 25 percent for 5-min averages. The shape of the spectrum computed from residuals is consistent with that derived from a model of solar plasma fluctuations. This data type can greatly aid navigation of a second spacecraft when the first is well-known in its orbit.

I. Introduction

The Venus-relative positions of the Magellan (MGN) and Pioneer 12 (PVO) spacecraft are currently being determined independently by using Earth-based measurements of the radio signals emitted by each spacecraft. Simultaneous tracking of the two orbiters provides much

stronger positioning information than is obtained by tracking just a single orbiter. The two spacecraft are so close angularly as seen from Earth that they may be observed in the same beamwidth of an Earth-based radio antenna. Radio metric data received at two Earth stations may be combined to provide an interferometric measurement of the plane-of-sky position difference between the two space-

craft. This measurement technique, called same-beam interferometry (SBI), is extremely precise due to double differencing of common errors along the four ray paths from the two spacecraft to the two ground stations. SBI measurements made at S-band (2.3 GHz) give plane-of-sky position change with an accuracy on the order of 10 m at Venus. A joint solution for the orbits of both spacecraft, combining Doppler and interferometric observables, provides up to an order of magnitude better accuracy than solutions for a single orbiter using only Doppler data.

A demonstration using the Magellan and Pioneer 12 orbiters is in progress to show the improvements to orbital accuracy that are provided by SBI, as compared with single-station Doppler. The first simultaneous data were acquired in August 1990, shortly after Magellan was inserted into Venus orbit, during the Magellan orbital check-out phase. SBI data were also acquired in February and April 1991. This article describes the data acquisition and presents an analysis of measurement system errors. An analysis of orbit solutions derived from these data will be presented in a following article.

SBI has been employed before, at the very beginning of the Pioneer 12 mission to Venus. Radio signals emitted from the four probes released by the orbiter were measured by SBI relative to the orbiter radio signal to determine the Venus wind speed and direction [1]. A similar determination was obtained from simultaneous interferometric measurements of the Vega spacecraft and the balloons that it dropped off at Venus [2]. Earlier, this technique was used to locate the Apollo 16 Lunar Rover relative to the Lunar Module [3] and to measure librations of the moon by using signals from the Apollo Lunar Surface Experiments Package (ALSEP) transmitters [4]. While these applications of SBI were of short duration and primarily for scientific purposes, NASA routinely and continuously tracks planetary orbiters. Obtaining telemetry and radio metric data simultaneously from two spacecraft at one station offers tremendous efficiency advantages. If, in addition, the two spacecraft are also simultaneously observed at a second, distant station during the mutual visibility period, then improved orbital accuracy results. Improved accuracy may increase science data return or reduce the total tracking time necessary to maintain a specified level of positional knowledge. The demonstration with Magellan and Pioneer 12 is the first step toward developing an operational capability for SBI data acquisition within NASA's Deep Space Network.

SBI system errors and the application of SBI to the positioning of planetary rovers, landers, and orbiters have been described previously [5-8]. When Doppler data are

supplemented with SBI data, the expected improvement in orbital accuracy for two orbiters varies from a factor of two to a factor of ten, depending on data accuracy and strategy. The Space Exploration Initiative (SEI) proposes a series of missions to Mars, which will likely involve communications orbiters, mapping spacecraft in low orbits, rendezvous craft, stationary landers, and rovers. Early development of an SBI capability will enhance this mission set.

The interferometric measurement and its expected contribution to the orbit solution process are discussed in more detail in the next section. Analytic models for predicting SBI measurement errors are presented. The acquisition and processing of SBI data from the Magellan and Pioneer 12 orbiters are briefly discussed. SBI measurement residuals from the August 1990 and February 1991 data sets are then analyzed and compared with predicted errors.

II. Tracking Strategies and SBI

Doppler measures the line-of-sight range rate between a Deep Space Station and a spacecraft. An arc of Doppler data from a planetary orbiter provides a history of the change in range. Interferometric measurements naturally complement Doppler. The difference in arrival time of a signal at two stations provides a direct measure of the angle between the baseline vector joining the two stations and the direction to the radio source. This, in turn, provides a geometric measure of the plane-of-sky position of the radio source in the direction of the baseline projected onto the plane of the sky.

The Doppler observable is derived from measuring the phase of the spacecraft RF carrier signal. Both Magellan and Pioneer 12 are capable of transmitting at S-band (2.3 GHz) and X-band (8.4 GHz). Since phase can be measured to within a small fraction of a cycle, Doppler is precise enough to sense submillimeter-level changes in line-of-sight range. It is generally not possible to utilize the full precision of Doppler since neither media fluctuations nor spacecraft dynamics can be modeled to the millimeter level. The SBI observable (Fig. 1) is also derived from measuring the carrier phase. The SBI observable is the phase, first differenced between stations and then differenced between spacecraft. Errors in the line-of-sight phase measurements are greatly diminished by double differencing. The doubly differenced phase, $\Delta^2\phi$, is given approximately by

$$\lambda \Delta^2\phi = (B \sin \theta) \Delta\theta$$

where

λ = the signal wavelength

B = the length of the baseline

θ = the angle between the baseline and the direction to the spacecraft

$\Delta\theta$ = the difference in θ for the two spacecraft

The angular separation between spacecraft is measured with an accuracy given by the doubly differenced phase accuracy times the wavelength divided by the baseline length projected onto the plane of the sky. The plane-of-sky position accuracy (linear distance at the spacecraft) of an SBI measurement is equal to the Earth-spacecraft distance times the angular accuracy. For millicycle-phase accuracy, intercontinental baselines, and an Earth-spacecraft distance of one astronomical unit, the position accuracy is on the order of 1 m. The combination of complementary Doppler and SBI data allows for dramatic improvement in orbital accuracy as compared with Doppler alone. Actual improvement in accuracy will be sensitive to where in the orbit the data are acquired.

Because the SBI observable is obtained from the carrier phase, there is an integer cycle ambiguity in the SBI measurement. A priori knowledge of orbiter positions will generally be inadequate to resolve this ambiguity. Thus, SBI measures the change in plane-of-sky position, rather than absolute position, just as Doppler measures the change in line-of-sight position. An SBI bias, common to all data points in a continuous arc, must be estimated. If the error in the estimation of the bias is significantly less than one cycle, it may be possible to fix the integer cycle, thus strengthening the data set. All instrumental effects not common to both spacecraft signal chains must be calibrated if the bias is to be fixed.

The orbits of both spacecraft are estimated together when using SBI data. The Pioneer 12 orbit is less perturbed by Venus gravity-field mismodeling than is the short-period Magellan orbit. For example, given a good orbit for Pioneer 12, SBI data will geometrically transfer position accuracy to Magellan, independent of gravity mismodeling. On the other hand, if Magellan is tracked much more frequently than Pioneer 12, then the SBI data improve the Pioneer 12 orbit by tying it to the relatively well-known Magellan orbit. In general, the signature in SBI data will provide information to improve the orbits of both spacecraft.

Currently, for operational tracking of either Magellan or Pioneer 12, one Deep Space Station is in two-way com-

munication with one spacecraft. The station transmits a signal that is coherently transponded by the spacecraft and received back at the transmitting station. The resulting measurement is referred to as two-way Doppler; orbit determination relies primarily upon this measurement. Orbit solutions for Magellan also make use of Doppler data received simultaneously at two stations and then differenced [9]. This data type, known either as narrowband VLBI [10], or more recently as differenced Doppler [11], has information content similar to that given by SBI, but is less accurate since measurement errors are not reduced by double differencing. For the SBI demonstration, Doppler data were acquired in the two-way mode from both Magellan and Pioneer 12, at separate stations. SBI data were acquired in an open-loop mode by using the Narrow Channel Bandwidth (NCB) VLBI System [12]. The NCB recordings were normally made on a noninterference basis at the Deep Space Stations scheduled for operational Magellan tracks. For some passes a separate station was scheduled for the purpose of recording with the NCB system. SBI data were acquired while both spacecraft were transmitting in the two-way mode, except for two passes when Magellan data were acquired in the one-way mode. The spacecraft onboard oscillator is the frequency reference for one-way Doppler and SBI. Because the received spacecraft signals are differenced between two stations, the SBI data accuracy does not depend on whether the tracking mode is one-way or two-way.

To realize advantages in efficiency by using SBI, it is of paramount importance to receive data from two or more spacecraft at a single station, with a second station to be used only during the baseline overlap periods. This mode of operation calls for replacing some or all of the two-way Doppler used for orbit determination with one-way Doppler so that multiple uplinks are not required, and it calls for simultaneous reception of telemetry from two spacecraft at one station. This is currently impractical for two reasons: (1) The oscillators on board the Magellan and Pioneer 12 spacecraft have insufficient stability to make one-way Doppler a viable alternative to two-way Doppler¹ and (2) the Deep Space Stations are configured to process only one telemetry signal stream at S-band and only one telemetry signal stream at X-band. In fact, dual support passes have been scheduled during which a single station receives Pioneer 12 telemetry at S-band and Magellan telemetry at X-band, but Doppler data are gen-

¹ D. Engelhardt, "Suitability of the MGN Onboard Auxiliary Oscillator for One-Way Doppler Data Generation," JPL Interoffice Memorandum 314.6-695 (internal document), Jet Propulsion Laboratory, Pasadena, California, July 28, 1986.

erated for only one spacecraft.² Future developments are expected to allow the DSN to fully realize the increased efficiency which SBI makes possible. Improvements to the stability of flight oscillators may enable orbit determination accuracy requirements to be met by using a combination of one-way Doppler and SBI data in place of two-way Doppler data [13]. The implementation of the Block V receiver in the DSN will enable simultaneous reception of Doppler and telemetry from at least two spacecraft at a single station at either S-band or X-band.³

III. Predicted SBI Measurement Errors

Measurement errors are predicted here for the SBI data obtained from MGN and PVO. Due to limited spacecraft battery capacity during data acquisition periods, Pioneer 12 was transmitting at S-band only, while Magellan was usually transmitting at both S-band and X-band. SBI data could be generated only at the S-band frequency, but the dual frequency downlink on Magellan was useful for characterizing charged particle delays. Measurements were compressed to either 2-sec, 20-sec, or 5-min averages over the 1-hr (typical) data arcs. Formulae for error prediction are developed that apply to arbitrary geometries and arbitrary frequencies. Solar plasma and instrumental phase shifts within the receiver are the two errors that are least perfectly canceled by double differencing, and generally they are limiting error sources for SBI measurements. These two error sources will be discussed in detail. For S-band measurements, the Earth's ionosphere is also a potentially significant error.

Distinction is made between (1) noiselike errors, which are random over a data arc and decrease with averaging time, and (2) errors that appear as slow drifts over a data arc. For this experiment, a constant measurement bias does not contribute an error, since an SBI data bias will be estimated for each data arc. Error predictions for SBI observables are given in picoseconds. One picosecond of time delay corresponds to 0.0023 cycle of phase at S-band, or to 0.0084 cycle of phase at X-band. Also, one picosecond of time delay corresponds to an angle of 37.5 prad over an 8000-km baseline, or to a transverse position of 8 m at 1.4 AU. Measurement errors are predicted for S-band SBI

observables for the geometries of August 1990 and February 1991. For comparison, error predictions are also given for X-band SBI observables for the geometry of August 1990. The assumptions used in the error analysis are summarized in Tables 1-4. A histogram of error contributions is shown in Fig. 2. All errors are one-sigma.

At S-band, the signal measured for SBI is the carrier from PVO and the carrier signal transmitted from either the high-gain antenna (HGA) or medium-gain antenna from MGN. The S-band carrier signals are separated by about 4 MHz in frequency. At X-band, the signal measured for SBI is the carrier from PVO and the -16th harmonic of the high-rate telemetry subcarrier from Magellan. The -16th harmonic is chosen as the Magellan X-band reference signal because it has adequate power and is separated by only 0.1 MHz in frequency from the PVO X-band carrier, whereas the Magellan X-band carrier is offset by 15 MHz from the PVO X-band carrier. Magellan was transmitting at S-band from its medium-gain antenna on August 11-12, 1990, and was transmitting dual-frequency S- and X-band downlinks from its HGA in February and April 1991.

A. System Thermal Noise

The sampled radio band will contain both the spacecraft signal and also the ground-receiver generated noise, which is proportional to the system operating temperature. The system noise error depends on the ratio of received signal power to system noise power. The voltage signal-to-noise ratio (SNR) for one-bit sampling is given by

$$SNR_V = \sqrt{\frac{4 P_{Tone}}{\pi N_0} T_{obs}}$$

where P_{Tone} is the received spacecraft tone power, N_0 is the system noise power in a 1-Hz bandwidth, and T_{obs} is the averaging interval (in seconds). The SBI thermal noise error is given by

$$\epsilon_{SBI} = 10^{12} \sqrt{2} \frac{1}{2\pi} \sqrt{\frac{1}{SNR_{V1}^2 f_{RF1}^2} + \frac{1}{SNR_{V2}^2 f_{RF2}^2}} \text{ psec}$$

where SNR_{Vi} and f_{RFi} are, respectively, the SNR_V and radio frequency in hertz for spacecraft i . The leading $\sqrt{2}$ accounts for two stations. A link analysis showing the received signal power and resulting voltage SNR for 1-sec integration is given in Table 4. Recorder data are time-multiplexed between spacecraft, so that 1 sec of data is

² *Deep Space Network Test and Training Plan, Magellan/Pioneer-12 Simultaneous Support (MPSS)*, JPL 870-176 (internal document), Jet Propulsion Laboratory, Pasadena, California, August 20, 1990.

³ *Deep Space Communications Complex, Block V Receiver Implementation Task, Volume 5: Functional Requirements*, JPL 803-112 (internal document), Jet Propulsion Laboratory, Pasadena, California, May 1, 1991.

recorded for each spacecraft in a 2-sec interval. The SNR_V for 20-sec averages is $\sqrt{10} \times SNR_{V,1sec}$, and the SNR_V for 5-min averages is $\sqrt{150} \times SNR_{V,1sec}$. System noise is a random error source. Though not decreased by station differencing or spacecraft differencing, system noise is reduced as averaging time increases.

B. Instrumental Phase Shifts

The instrumental phase response of the receiver system is characterized predominantly as a group delay: The instrumental phase shift is a linear function of the frequency of the received signal. Deviations from linearity due to dispersion take two forms: a curvature that slowly varies with signal frequency and a much more rapidly varying "ripple." SBI measurements are affected in a systematic way as the Doppler shift causes the received signal frequency to sweep across the passband. Phase shifts due to instrumental group delays and clock offsets are referred to as linear systematic; phase shifts due to bandpass curvature are referred to as nonlinear systematic; and shifts due to the deviation of the system phase response from a smooth phase response are referred to as phase ripple.

1. **Linear Systematic.** The time-dependent station-differenced phase shift $\phi_I(t)$ for a single spacecraft due to instrumental group delays and station clock epoch offsets can be shown to be⁴

$$\phi_I(t) = f_{RF}(\nu(t)\tau_I + \dot{\rho}(t)\delta\tau_{clock} - \delta\tau_{clock}) + \Delta\phi_h \text{ cycles}$$

where f_{RF} is the transmitter frequency (in hertz); $\nu(t)$ is the interferometer fringe frequency (in sec/sec); $\dot{\rho}(t)$ is the line-of-sight range rate to Station 1 (in sec/sec); τ_I is the instrumental group delay at Station 2 (in seconds); $\delta\tau_{clock}$ is the error in modeling the clock epoch offset between stations (in seconds); and $\Delta\phi_h$ is a constant instrumental phase shift (in cycles). Note that $\nu(t)$ is the line-of-sight range rate differenced between stations. The term $\dot{\rho}(t)$ refers to Station 1, and the term τ_I refers to Station 2, due to the convention that was used to select a clock reference point. The phase shift $\phi_I(t)$ causes an SBI error because the model for spacecraft-received phase does not account for certain instrumental delays and, of course, for an unknown clock offset. A constant phase shift is absorbed by estimating an SBI data bias. SBI errors will be induced by a change in ϕ_I over the data arc, which is not common to both sources. The change in ϕ_I over a data arc due to the change in ν and $\dot{\rho}$ is

⁴ J. S. Border, "Analysis of Clock and Instrumental Group Delays in Δ V LBI Observables," JPL Interoffice Memorandum 335.1-90-002 (internal document), Jet Propulsion Laboratory, Pasadena, California, January 5, 1990.

$$\delta\phi_I = f_{RF}(\delta\nu\tau_I + \delta\dot{\rho}\delta\tau_{clock}) \quad (1)$$

The interferometer fringe frequency ν , approximately given by the Earth rotation rate times the baseline length divided by the speed of light, is almost independent of spacecraft velocity for a distant spacecraft. The change in ν over a pass may be $\delta\nu = 0.3 \mu\text{sec}/\text{sec}$. The instrumental group delay τ_I is about $25 \mu\text{sec}$. But $\delta\nu$ is almost the same for two spacecraft at Venus, the difference being on the order of $1 \text{ nsec}/\text{sec}$. This makes the first term in Eq. (1) an insignificant SBI error. The change in line-of-sight range rate $\delta\dot{\rho}$ will be distinct for each spacecraft, unless they share an orbit (e.g., lander and rover). For SBI data arcs acquired in August 1990, the maximum change in $\delta\dot{\rho}$ occurred for Magellan and was $22 \mu\text{sec}/\text{sec}$. A clocklike offset between stations can be measured with subnanosecond precision by making interferometric observations of natural radio sources. However, the term $\delta\tau_{clock}$ above refers to the error in knowledge of time tags. Currently this knowledge is on the order of $0.1 \mu\text{sec}$. Ambiguities in instrumental calibrations have prevented absolute time-tag synchronization at the precision of the measurements. This is not a fundamental limitation; nanosecond-level synchronization is entirely feasible [14]. The second term in Eq. (1) is one of the largest SBI errors; improved synchronization is, in fact, required to reduce this error source. The SBI error due to linear systematic effects is given by

$$\epsilon_{SBI} = 10^{12} \delta\dot{\rho} \delta\tau_{clock} \text{ psec}$$

This error changes slowly over a pass, as $\dot{\rho}$ changes.

2. **Nonlinear Systematic.** The phase response of the baseband filters has a known systematic curvature. The curvature is almost the same at each station. The received baseband frequencies will, in general, differ at each station due to the rotation of the Earth, so the station-differenced phase shift will not be zero. The phase shift will change as the narrowband signal sweeps across the passband. To estimate the size of this effect, suppose that the curvature is quadratic over $2F$ kHz with a peak-to-peak phase nonlinearity of $2A$ deg. Suppose that the baseband frequencies at the two stations are separated by δf_{bb} kHz. Then, the maximum station-differenced phase shift as the received signal sweeps from $-F$ to $+F$ kHz at Station 1 is

$$\delta\phi_c \approx \pm \frac{4A\delta f_{bb}}{F} \text{ deg}$$

In August 1990, the Magellan received frequency changed by about 50 kHz over a 1-hr data arc, while the

PVO received frequency changed by less than 3 kHz. The peak-to-peak nonlinear phase response over 50 kHz is estimated to be 2 deg for the open-loop receiver system used [15]. The offset in baseband frequencies between stations was 4 kHz for data acquired in August 1990. This gives a phase shift of $\delta\phi_c = 0.64$ deg.

The SBI error is given by

$$\epsilon_{SBI} = 10^{12} \frac{\delta\phi_c / 360 \text{ deg}}{f_{RF}} \text{ psec}$$

This error varies slowly, as the line-of-sight range rate varies, and does not cancel between spacecraft. The error can be made smaller by either modeling the curvature of the passband or by offsetting station mixing frequencies so that baseband frequencies will be nearly equal. For data acquisitions in February and April 1991, station mixing frequencies were offset, virtually eliminating the error due to bandpass curvature at baseband. There also is an error of this form due to curvature at RF and IF, which is not eliminated by offsetting station mixing frequencies, though the magnitude of this curvature is estimated to be smaller than the curvature at baseband by a factor of ten.

3. Phase Ripple. The phase ripple ϵ_ϕ for the open-loop system used to record the SBI data is estimated to be 0.5 deg [15]. Most of the ripple comes from the baseband filter, though RF and IF components also contribute. Variations of ± 0.5 deg away from the smooth systematic curvature of the bandpass occur over scales of a few kilohertz. Thus, this error is generally independent for each spacecraft and each station, since received frequencies are not equal. The SBI error is given by

$$\epsilon_{SBI} = 10^{12} \frac{2\epsilon_\phi / 360 \text{ deg}}{f_{RF}} \text{ psec}$$

Received frequencies from Venus orbiters may change by a few kilohertz in 5 min, or by much less, so this error is random or systematic, depending on the orbital geometry. For the MGN-PVO data set of August 1990, the error component due to Magellan is random, while the error component due to PVO is slowly varying. For a data arc that includes the time of PVO periapsis, the error component due to PVO is random rather than systematic.

C. Transmitter Frequency Offset

Clock offsets and instabilities at both transmitters and receivers are nearly eliminated by differencing. But the

transmitter frequency is used to convert between calculated geometric time delays and observed phase delays. An unknown offset δf_T (in hertz) between the transmitter frequencies of two spacecraft will cause an SBI error of

$$\epsilon_{SBI} = 10^{12} \tau \frac{\delta f_T}{f_{RF}} \text{ psec}$$

where τ (in seconds) is the interferometric delay. For two-way signals, where the spacecraft is transponding a signal uplinked from a ground station, the offset $\delta f_T / f_{RF}$ is a measure of clock-rate synchronization between the two stations that are uplinking to the two spacecraft. The DSN station clock rates are generally synchronized to 10^{-13} sec/sec, though a few stations have independent clocks that may have unknown rate offsets as large as 10^{-12} sec/sec. The interferometric delay may range over -0.02 to $+0.02$ sec. This error source is insignificant for two-way transmissions, such as were used for nearly all the Magellan and Pioneer 12 measurements. For one-way transmissions, the spacecraft onboard oscillator frequency must be estimated to the 0.1-Hz level in order to reduce the size of this error below one picosecond.

D. Baseline

Baseline errors include uncertainties in station location, Earth orientation, and frame tie. The DSN Earth-fixed baseline components are known to 5 cm. Calibrations for UT1-UTC, polar motion, and nutation tweaks are available about 2 weeks after the measurement epoch that provides Earth orientation accuracy of 5 cm per component in the radio frame. Real-time predictions are currently at the 30-cm level, but may improve to 5 cm over the next few years as Global Positioning System (GPS) satellite measurements are incorporated into Earth orientation solutions [16,17]. Since the positions of the spacecraft being tracked are specified relative to Venus, the orientations of the baselines in the planetary frame are required. Knowledge of the orientation offset between the radio frame and the planetary frame is estimated to be 25 nrad [18]. This causes a baseline error of 20 cm per component. The rss baseline error (for data processed 2 weeks after real time) is 21.2 cm per component. The SBI error is given by

$$\epsilon_{SBI} = \frac{10^{12}}{c} \epsilon_{BL} \Delta\theta \text{ psec}$$

where ϵ_{BL} is the baseline component error (in centimeters), $\Delta\theta$ is the angular separation of the two spacecraft (in radians), and c is the speed of light (cm/sec). This error term is slowly varying over a pass.

E. Troposphere

Signal path delay through the troposphere is a function of the zenith tropospheric delay and source elevation angle. The zenith delay is computed from surface meteorology data, mapped to the line of sight, and applied as a calibration for each source. The calibration error depends on the uncertainty in the determination of zenith delay and the differential elevation angle between sources. Uncertainty in the static component of the zenith delay is dominated by the uncertainty in determining the wet component. Spatial fluctuations in water vapor content also affect SBI observables.

1. **Zenith Bias.** At a single station, the error due to a zenith bias is given by

$$\epsilon_{SBI} = \frac{10^{12}}{c} \frac{\rho_{zi} \cos \gamma}{\sin^2 \gamma} \delta \gamma \text{ psec}$$

where ρ_{zi} is the uncertainty in the measurement of zenith delay (in centimeters), c is the speed of light (cm/sec), γ is the source elevation angle, and $\delta \gamma$ is the differential elevation angle between sources (in radians). This error will be independent at each station and will vary slowly over a pass. It is calculated for the MGN-PVO data set from the parameters in Tables 2 and 3.

2. **Fluctuations.** Two ray paths from a single station, separated by $\Delta \theta$ rad, may be thought of as having a spatial separation of $3h_t \Delta \theta$ in the tropospheric shell, where h_t is the effective wet tropospheric height ($h_t \approx 1$ km) and 3 is the mapping to the typical elevation angle of 20 deg. An estimate for the effect of spatial fluctuations on two ray paths separated by this distance may be derived from the structure function of tropospheric delay developed by Treuhaft and Lanyi [19]. Fluctuations are described by Kolmogorov turbulence. For small angular separations ($\Delta \theta \leq 1$ mrad), expressed as a function of the angle $\Delta \theta$ (in radians), the SBI error is given by

$$\epsilon_{SBI} = \sqrt{2} 48 \Delta \theta^{5/6} \text{ psec}$$

$$\text{decorrelation time} = (3h_t \Delta \theta) / v_t \text{ sec}$$

where $v_t \approx 0.008$ km/sec is the tropospheric wind speed and the leading $\sqrt{2}$ accounts for two stations. For the August 1990 MGN-PVO data set, with $\Delta \theta = 0.3$ mrad, the SBI error is 0.08 psec with a decorrelation time of 0.11 sec. This error is insignificant.

F. Ionosphere

The ionosphere is a dispersive medium; path delay scales as the inverse of frequency squared. Charged particle effects can be eliminated by observing at two radio frequencies, such as S- and X-band. For single-frequency data, calibrations are applied based on measurements of the total electron content (TEC) along the line of sight from each Deep Space Station to one or more beacon satellites. A mapping function, which depends on elevation angle and on the arc length in the ionospheric shell between the Earth-sun line and the source ray path, is used to calculate calibrations for each source. Systematic errors are caused by uncertainty in the measured TEC to the beacon satellite, by mapping error, and by a difference in radio frequency between the two observed sources. Spatial fluctuations in TEC also affect SBI observables.

1. **Zenith Bias.** For a single source, the total ionospheric delay can be written as

$$\tau_{ION} = 1340 \frac{f TEC_z}{f_{RF}^2} \text{ psec}$$

where TEC_z is the zenith delay in TEC units (10^{16} el/m²), f is the mapping function from zenith to line of sight, and f_{RF} is the radio frequency in gigahertz. The differential error, after calibration, for two sources observed at one station is then given by

$$\epsilon_{SBI} = 1340 \frac{\rho_{zi}}{f_{RF}^2} \left(\frac{\partial f}{\partial \theta} \Delta \theta + 2f \frac{\delta f_{RF}}{f_{RF}} \right) \text{ psec} \quad (2)$$

where ρ_{zi} is the error in the zenith delay measurement (in TEC units), $\partial f / \partial \theta$ is the spatial derivative⁵ of the mapping function, and δf_{RF} (in gigahertz) is the difference in frequency between the two observed sources. Zenith delay calibrations were obtained from Faraday rotation measurements of geostationary satellite beacons and from dual-frequency group delay measurements of the available GPS satellites. The zenith error ρ_{zi} is estimated to be 5 TEC units. To estimate $\partial f / \partial \theta$ for the MGN-PVO data set, calibrations mapped to line of sight were examined for the full Venus visibility window for each station complex. The maximum change in the mapping function per unit change in angle was found to be $\partial f / \partial \theta = 6.9 \text{ rad}^{-1}$. For

⁵ The mapping function f depends on the two-dimensional angle that defines the ray path, and it depends on diurnal variations in zenith TEC. Here, the zenith error ρ_{zi} is considered to be constant and $\partial f / \partial \theta$ is computed as the maximum (over all directions) change in mapped calibration per unit angle.

$\Delta\theta = 0.3$ mrad, the first term in Eq. (2) is 2.6 psec. This term slowly varies over a pass, as the zenith measurement error and the geometry change. For a typical mapping of $f = 3$, the second term in Eq. (2) is 13.2 psec. But this is the total effect; a constant offset is absorbed by estimation of an SBI data bias. The SBI data are only sensitive to a change in the delay due to this term. The mapped calibration changed by at most 25 TEC units over a MGN-PVO pass. The change in the calibration error is expected to be 10 percent of this. Substituting 2.5 TEC units for $f\rho_{zi}$ in the second term gives an SBI error of 2.2 psec. The total systematic error will not be increased by $\sqrt{2}$ since neither the spatial derivative of the mapping function nor the drift in calibration error is expected to be large at both complexes concurrently.

Ionospheric calibration uncertainty is a potentially significant error source for S-band SBI measurements. This source of error is expected to be reduced by a factor of 2 to 5 when dual-frequency GPS group-delay measurements become available from a full GPS constellation for generation of ionospheric calibrations. Also, the component of error dependent on frequency difference between sources can be reduced if even one source transmits dual-frequency S-X signals. The dual-frequency transmissions allow precise measure of the change in TEC along the line of sight. The Magellan spacecraft transmitted dual-frequency S-X signals during the February and April 1991 SBI measurement sessions. The dual-frequency data were used to check the externally supplied ionosphere calibrations. For this data set, the observables were changed by an insignificant amount when the Faraday/GPS calibrations were adjusted so that the change in TEC along the line of sight was as measured by the Magellan S-X data.

2. Fluctuations. Temporal fluctuations in the ionospheric delay rate, after removal of a nominal calibration, have been observed to be about 3.8×10^{-14} sec/sec for X-band signals, for daytime low-elevation observations, and for averaging intervals from 60 to 6000 sec.⁶ An estimate for SBI errors due to spatial fluctuations will be derived from this result. The angular separation $\Delta\theta$ (in radians) between ray paths corresponds to a spatial separation d (in kilometers) in the ionospheric shell and to a temporal separation T (in seconds) between measurements through the relations

$$d \approx Tv_i \approx 3h_i\Delta\theta$$

where $v_i \approx 0.1$ km/sec is the ionospheric wind speed, 3 is a typical elevation scale factor, and $h_i \approx 350$ km is the effective height of the ionosphere. Spatial fluctuations in delay, for angles from 0.006 to 0.6 rad, appear to scale linearly with the angle, since delay rate variations are flat over corresponding time intervals. It may be optimistic to assume that the angular dependence remains linear for very small angles. Instead, Kolmogorov turbulence will be assumed. Thus, the error is of the form $k\Delta\theta^{5/6}$ for small angles. The constant k is defined by making this expression predict an error for $\Delta\theta = 0.00571$ rad, which is consistent with the corresponding observed temporal variations for a time offset of 60 sec. The resulting SBI error is

$$\epsilon_{SBI} = \sqrt{2} \frac{11700}{f_{RF}^2} \Delta\theta^{5/6} \text{ psec}$$

$$\text{decorrelation time} = (3h_i\Delta\theta)/v_i \text{ sec}$$

where f_{RF} is the observing frequency in gigahertz and $\sqrt{2}$ accounts for two stations. For the August 1990 MGN-PVO data set, with $\Delta\theta = 0.3$ mrad, the SBI error is 3.6 psec with a decorrelation time of 3.15 sec.

G. Solar Plasma

To model solar plasma-induced fluctuations on radio signals transmitted by interplanetary spacecraft, the plasma is imagined to be confined to a thin screen passing through the center of the sun and perpendicular to the line connecting the spacecraft and the Earth. Define a random function, $\phi(x)$, which represents the phase fluctuation induced on a radio signal as it penetrates the plasma screen at a distance x from the center of the sun. Then the quantity $\phi(b+x) - \phi(x)$ represents the differential phase fluctuation induced on two signals that are separated by a distance b when passing through the plasma screen. For an interferometric observation, b is the length of the projection of the baseline onto the plasma screen. For an observation of a spacecraft at Venus from a DSN baseline, b is typically 4000 km.

To determine the plasma-induced phase error on interferometric measurements, the spatial structure function of phase, $D(b)$, is computed. It is defined as: $D(b) \equiv \langle (\phi(b+x) - \phi(x))^2 \rangle$. Here, brackets denote an ensemble average. The value $D(b)$ may be viewed as representing the phase variance of a station-differenced phase measurement. The power spectrum of electron density fluctuations, which has been determined experimentally [20], can be used to calculate a formula for $D(b)$ [21]

⁶ A. J. Mannucci, "Temporal Statistics of the Ionosphere," JPL Interoffice Memorandum 335.1-90-056 (internal document), Jet Propulsion Laboratory, Pasadena, California, October 25, 1990.

$$D(b) = \frac{2.5 \times 10^{-4}}{(f_{RF})^2} \times (b/v_{sw})^{1.65} \times (\sin(SEP))^{-2.45} \text{ cycles}^2 \quad (3)$$

In this expression, f_{RF} is the signal radio frequency in gigahertz, SEP is the sun-Earth-probe angle, and v_{sw} is the velocity of the solar wind (typically 400 km/sec).

Knowledge of the structure function enables calculation of the temporal correlations between plasma-induced errors on interferometric measurements. First note that because of the dynamics of the solar wind, the plasma-induced phase fluctuation, ϕ , is actually a function of both space and time, $\phi = \bar{\phi}(x, t)$. If, however, it is assumed that the plasma turbulence consists of fixed structures that maintain their shape as they travel radially outward from the sun at velocity v_{sw} , then ϕ may be written as a function of a single variable $\bar{\phi}(x, t) = \phi(x - v_{sw}t)$.

Let $\Delta\Phi(b, T_i) \equiv \bar{\phi}(b+x, t_i) - \bar{\phi}(x, t_i)$ and $\Delta\Phi(b, T_j) \equiv \bar{\phi}(b+x, t_j) - \bar{\phi}(x, t_j)$ be the station-differenced phase at Earth receive times T_i and T_j , along a baseline whose projection onto the plasma screen has length b . Here, $t_i = T_i - [\text{light travel time from plasma screen to Earth}]$, and $t_j = T_j - [\text{light travel time from plasma screen to Earth}]$. (The slight difference in signal arrival time at the two stations is ignored here.)

Then, the temporal covariance is computed as

$$\begin{aligned} \langle \Delta\Phi(b, T_i) \Delta\Phi(b, T_j) \rangle &= \langle (\bar{\phi}(b+x, t_i) - \bar{\phi}(x, t_i)) \\ &\quad \times (\bar{\phi}(b+x, t_j) - \bar{\phi}(x, t_j)) \rangle \\ &= \langle \bar{\phi}(b+x, t_i) \bar{\phi}(b+x, t_j) \rangle \\ &\quad - \langle \bar{\phi}(b+x, t_i) \bar{\phi}(x, t_j) \rangle \\ &\quad - \langle \bar{\phi}(b+x, t_j) \bar{\phi}(x, t_i) \rangle \\ &\quad + \langle \bar{\phi}(x, t_i) \bar{\phi}(x, t_j) \rangle \end{aligned}$$

Assuming that ϕ is stationary,

$$\begin{aligned} \langle \bar{\phi}(x, t) \bar{\phi}(y, t') \rangle &= \langle \phi(x - v_{sw}t) \phi(y - v_{sw}t') \rangle \\ &= \langle \phi(x - y - v_{sw}(t - t')) \phi(0) \rangle \end{aligned}$$

and

$$D(\rho) = \langle (\phi(\rho+x) - \phi(x))^2 \rangle = 2 \{ \langle \phi^2 \rangle - \langle \phi(\rho)\phi(0) \rangle \}$$

Applying the above two equations yields

$$\begin{aligned} \langle \Delta\Phi(b, T_i) \Delta\Phi(b, T_j) \rangle &= \langle \phi^2 \rangle - 1/2D(v_{sw}(t_j - t_i)) \\ &\quad - \langle \phi^2 \rangle + 1/2D(b - v_{sw}(t_j - t_i)) \\ &\quad - \langle \phi^2 \rangle + 1/2D(b + v_{sw}(t_j - t_i)) \\ &\quad + \langle \phi^2 \rangle - 1/2D(v_{sw}(t_j - t_i)) \\ &= 1/2D(b - v_{sw}(t_j - t_i)) \\ &\quad + 1/2D(b + v_{sw}(t_j - t_i)) \\ &\quad - D(v_{sw}(t_j - t_i)) \end{aligned}$$

Since D can be computed by using Eq. (3) above, this expression for the temporal covariance between station-differenced phase measurements can be used to estimate the plasma-induced error on an interferometric observable that is derived from a time average of "instantaneous" measurements. As an example, consider a spacecraft at Venus on August 11, 1990, transmitting an S-band signal that is simultaneously received at Goldstone and Canberra. The SEP angle is 21 deg, and the projection of the baseline onto the plasma screen is 4000 km. The above model for the solar plasma yields a plasma-induced station-differenced delay error of 40 psec over 5 min.

The SBI observable is formed from doubly differenced phase observables. For spacecraft whose separation is less than 1 deg, double differencing of phase observables can result in significant cancellation of the plasma-induced error. Let $\Delta^2\Phi(S, b, T_j)$ be the doubly differenced phase at time T_j . Here S is the separation (on the plasma screen) between the signals transmitted by the two spacecraft (Fig. 3). Then, the instantaneous phase variance is

$$\begin{aligned} \Delta^2\Phi(S, b, T_j) &\equiv \{ \Delta\Phi(b, T_j) \}_{S/C1} - \{ \Delta\Phi(b, T_j) \}_{S/C2} \\ &= \{ \bar{\phi}(S+b+x, t_i) - \bar{\phi}(S+x, t_i) \} \\ &\quad - \{ \bar{\phi}(b+x, t_i) - \bar{\phi}(x, t_i) \} \end{aligned}$$

The covariance between doubly differenced phase measurements, $\langle \Delta^2\Phi(S, b, T_i) \Delta^2\Phi(S, b, T_j) \rangle$, can be com-

puted by first applying the substitution above and then proceeding in an analogous fashion to the calculation

for singly differenced phase measurements. The final result is

$$\begin{aligned} \langle \Delta^2 \Phi(S, b, T_i) \Delta^2 \Phi(S, b, T_j) \rangle = & -2D(v_{sw}(t_j - t_i)) + D(b + v_{sw}(t_j - t_i)) + D(S + v_{sw}(t_j - t_i)) \\ & + D(-b + v_{sw}(t_j - t_i)) + D(-S + v_{sw}(t_j - t_i)) \\ & - 1/2D(b + S + v_{sw}(t_j - t_i)) - 1/2D(S - b + v_{sw}(t_j - t_i)) \\ & - 1/2D(b - S + v_{sw}(t_j - t_i)) - 1/2D(-b - S + v_{sw}(t_j - t_i)) \end{aligned}$$

D can be computed by using Eq. (3) above, so this expression for the temporal covariance between doubly differenced phase measurements can be used to estimate the plasma-induced error on an SBI observation.

In the August 1990 SBI experiment, the SEP angle was 21.3 deg. The baseline projection onto the plasma screen, b , was about 3300 km for Goldstone-Madrid and 4000 km for Goldstone-Canberra; the projection of the spacecraft separation, S , was about 40,000 km during the Goldstone-Madrid observations and 30,000 km during the Goldstone-Canberra observations. Using the above model for the solar plasma, the plasma-induced SBI error for a 5-min average is 13 psec at S-band for each baseline. This is substantially smaller than the 40-psec plasma-induced delay error for a single spacecraft interferometric observable; because Magellan and PVO are angularly close (separation ≈ 0.3 mrad), plasma-induced phase advances on signals transmitted by each spacecraft are highly correlated. The predicted SBI error for a 2-sec average is 81 psec, while the error for a 20-sec average is 64 psec. Errors have high temporal correlation for averaging times less than 100 sec.

This calculation of the covariance between doubly differenced phase measurements assumes that the two spacecraft transmit signals at exactly the same frequency; Magellan and PVO transmit signals at 2297 GHz and 2293 GHz, respectively. Since (1) plasma-induced phase advance is inversely proportional to frequency; (2) the spacecraft frequencies differ by only 0.2 percent; and (3) the predicted doubly differenced phase error is only about a factor of three less than the singly differenced phase error, the temporal covariance of doubly differenced phase measurements differs from the above estimate by less than one percent.

It should be noted that the separations S and b typically have components that are not in the solar radial

direction. It is expected that the error for such a measurement is comparable with the error calculated via the simple model above. Large-scale turbulence ($>40,000$ km or, equivalently, hundreds of seconds) is still common to all four signal paths; small-scale turbulence (less than a few thousand kilometers or, equivalently, several seconds) is not common to any of the signal paths, regardless of the orientation of S and b . Thus, for averaging times of 2 sec or 5 min, the predicted plasma-induced error on the SBI measurement should not depend strongly on the orientation of S and b .

SBI observations of two spacecraft with yet smaller angular separation would benefit from even greater cancellation of the solar plasma error. SEI encompasses an extensive series of missions to Mars beginning in the late 1990s. Two spacecraft separated by thousands of kilometers on the Martian surface would have an angular separation on the order of 5–10 μ rad as viewed from Earth; when Mars is at a 20-deg SEP angle, the plasma-induced phase error for a 5-min S-band SBI observation of two landed spacecraft is less than 2.2 psec.

Solar plasma turbulence can vary substantially from day to day, potentially resulting in an order of magnitude variation in plasma-induced phase error at a given SEP angle [20].

H. Root-Sum-Square Error

The total predicted measurement error is computed as the rss of individual error terms. It is evident from Fig. 2 that solar plasma fluctuations dominate S-band SBI errors for the geometry of August 1990. This error term is reduced for the geometry of February 1991, due to the reduced angular separation between the two spacecraft as seen from Earth. Errors due to station instrumentation and the Earth's ionosphere are comparable to the solar

plasma error for the geometry of February 1991. The predicted one-sigma rss error for the S-band SBI measurements acquired in February 1991 from Magellan and PVO is, for averaging intervals of 2 sec, 20 sec, and 5 min, equal to 56 psec, 30 psec, and 4.8 psec, respectively. The plane-of-sky separation between the two spacecraft is thus measured with a precision of 41 m, after a 5-min integration. All observables in a continuous data arc will have a common bias, due to the integer cycle ambiguity in the measurement of radio signal phase.

IV. Data Acquisition and Processing

SBI data may be acquired during the overlap periods when Venus is visible from either the Goldstone-Madrid baseline or the Goldstone-Canberra baseline, provided that both Magellan and Pioneer 12 are transmitting. Neither spacecraft transmits continuously, however. Magellan, which has an orbital period of 3.25 hr, transmits to Earth for approximately 2 hr of each revolution. Data are unavailable when Magellan is making radar measurements of the Venus surface or performing a star calibration. Pioneer 12 transmission time is limited by the amount of power available at the spacecraft, which has been decreasing during the last few years. Daily transmission time has typically been in the range of 4 to 12 hr since August 1990. SBI data were scheduled when both Magellan and Pioneer 12 were expected to be transmitting during baseline overlaps.

Doppler data were acquired whenever possible. During each day, tracking shifts from the DSN complex at Madrid, Spain; to Goldstone, California; to Canberra, Australia. Reference orbits were generated separately for each spacecraft by using the two-way Doppler.

A total of 33.5 hr of SBI data were acquired from Magellan and Pioneer 12; the amount of data acquired during each experiment set for each baseline is listed in Table 5. Figure 4 shows the Venus-centered orbits of the two spacecraft for the geometry of August 1990, projected onto the plane of the sky, and indicates where in the orbits SBI data were acquired. Pioneer 12 was near apoapsis for all SBI data acquisitions in August 1990. During February and April 1991, Pioneer 12 was near periapsis during the Goldstone-Canberra overlap.

The NCB VLBI System was used to make an open-loop recording (1-bit samples) of the received signal voltage in 250-kHz channels that contained the S-band carrier signals from each spacecraft. The Magellan X-band carrier signal was also recorded in the February and April 1991 experiments. The signal phase was extracted at 1-sec intervals

by digitally mixing a model of the received phase with the recorded signal voltages. For each spacecraft, the phase was differenced between stations, then further compressed to either 20-sec or 5-min averages. Residuals were obtained by removing a model based on the reference trajectories and by applying calibrations for tropospheric and ionospheric delays. Tropospheric calibrations were based on surface meteorological data, and ionospheric calibrations were based on either Faraday rotation measurements or dual-frequency GPS satellite measurements. SBI residuals were then obtained by differencing between spacecraft. Using two-way Doppler and SBI, a simultaneous solution for the orbits of both spacecraft was then generated. For the purpose of examining measurement errors, residuals were recomputed with respect to the new spacecraft orbits.

V. Analysis of Measurement Residuals

The magnitude of measurement errors and the characterization of correlations among measurements are important inputs to the orbit determination process. SBI residuals are examined here in an attempt to validate the error budget. Since the spacecraft orbits were fit to the SBI data as well as to the longer Doppler data arcs, it is possible that systematic SBI measurement errors have been absorbed into the spacecraft orbits. Checks of orbit consistency and orbit prediction will be necessary to resolve this issue. Here, examination is restricted to residuals of individual data arcs.

Residuals of 2-sec averages are displayed in Fig. 5 for the longest contiguous data arc acquired in August 1990. The geocentric angular separation of the two spacecraft is also shown. The predicted one-sigma error is 82 psec for 2-sec averages. The rms of the residuals is less than the prediction by about a factor of three at this time scale. When data are further compressed to 5-min averages, the rms is less than the predicted error of 14 psec by about 25 percent.

Residuals of 20-sec averages are shown in Fig. 6 for SBI data acquired on the Goldstone-Canberra baseline on February 17, 1991. A dramatic drop in the point-to-point scatter occurs where the angular spacecraft separation (also shown in Fig. 6) reduces from 100 to 30 μ rad. This is believed to be the result of a more complete cancellation of solar plasma effects. Angular spacecraft separations below 100 μ rad did not occur during SBI data recordings in August 1990.

Figure 7 shows residuals of 20-sec averages for SBI data acquired on the Goldstone-Canberra baseline on February

20, 1991. The variation in angular separation over this data arc is comparable to the variation over the data arc displayed in Fig. 6, yet the apparent reduction in point-to-point scatter is much less dramatic at the point of closest angular approach. This may be due to temporal variations in solar wind intensity or to geometric factors not accounted for in the solar plasma error model developed here. But, for both the February 17 and the February 20 data arcs, the rms of residuals for 20-sec averages is less than that predicted by the error budget. The rms for February 17 is 29 psec; the rms for February 20 is 18 psec; while the error budget prediction (assuming 90- μ rad angular separation) is 30 psec.

According to the measurement error budget, solar plasma is the dominant error source at S-band frequencies for most geometries. Without dual-band measurements, it is difficult to separate dispersive and nondispersive errors in the residual phase; one method of determining whether the measurements are consistent with the error budget is to examine the power spectrum of the doubly differenced residual phase to see whether it conforms to the behavior expected of the solar plasma. Radio scattering experiments with the Pioneer and Viking spacecraft [20,22] have established that the one-dimensional power spectrum of solar plasma-induced phase scintillations follows a power law of the form $P_\phi(f)$ proportional to $f^{-2.65}$. Since the SBI observable is formed by doubly differencing four line-of-sight phase measurements, the form of the spectrum for the SBI observable will be somewhat modified; differencing the signals acts as a crude high-pass filter.

The function $\phi(x - v_{sw}t)$ has been defined in Section III.G as representing the plasma-induced phase advance on an individual signal traversing the plasma screen at a distance x from the sun at time t . Define $\varphi_x(t) = \phi(x - v_{sw}t)$. Then $\varphi_x(t)$ represents the temporal phase fluctuations induced by the plasma at a distance x from the sun. The power spectrum of the temporal phase fluctuations is the Fourier transform of the autocorrelation of φ_x . The power spectrum is of the form $Kf^{-2.65}$, where the value of the constant K depends on the distance x , or equivalently, the SEP angle. For a spacecraft at a 20-deg SEP angle transmitting a 2.3-GHz signal, $P_{\varphi_x}(f) = 3.9 \times 10^{-6} f^{-2.65}$ cycles²/Hz [20]. To represent a doubly differenced phase measurement, define a new function, $\xi_x(t) = \varphi_x(t + \tau_S + \tau_b) - \varphi_x(t + \tau_S) - \varphi_x(t + \tau_b) + \varphi_x(t)$, where τ_b is the length of the baseline projection on the plasma screen, b , divided by the solar wind velocity, and τ_S is the projection of the spacecraft separation on the plasma screen, S , divided by the solar wind velocity. The power spectrum of ξ_x can be derived by Fourier transform-

ing the autocorrelation of ξ_x . The result is given by the following expression:

$$P_{\xi_x}(f) = 2P_{\varphi_x}(f)[2 - 2\cos(2\pi\tau_b f) - 2\cos(2\pi\tau_S f) + \cos(2\pi(\tau_b + \tau_S)f) + \cos(2\pi(\tau_b - \tau_S)f)] \quad (4)$$

Note that when f is small, as compared with $1/\tau_b$ and $1/\tau_S$, a Taylor series expansion of the cosines in the above equation shows that $P_{\xi_x}(f)$ is proportional to $f^{+1.35}$. For higher frequencies, the power spectrum follows a power law of the form $f^{-2.65}$, which is modulated by cosines.

For the SBI experiment in August 1990, τ_b and τ_S were typically 10 sec and 75 sec (here a solar wind velocity of $v_{sw} = 400$ km/sec is assumed). Using these values, Eq. (4) is plotted in Fig. 8 with the power spectrum of SBI residuals from the Goldstone-Canberra data arc on August 11 [note that in this figure, the cosine minima that occur in Eq. (4) have been suppressed, and only the envelope of the theoretical spectrum is shown]. The experimental and theoretical spectra are qualitatively similar, though the experimental spectrum is not as strongly suppressed at low frequencies as is the theoretical spectrum. Other error sources in addition to solar plasma may be increasing the magnitude of the experimental spectrum at the lowest frequencies. Both spectra have a "knee" at close to 0.01 Hz (note that $v_{sw}/S = 0.013$ Hz). The theoretical spectrum is about one order of magnitude larger than the experimental spectrum for higher frequencies; this is consistent with the factor of 3 difference between the predicted phase error and scatter of the SBI residuals at short averaging intervals. This discrepancy could reflect either daily variation in the solar wind turbulence or it may suggest that the model for plasma-induced doubly differenced phase scintillations needs to be renormalized. The acquisition of additional data, particularly dual-band data, is needed to help resolve this issue. Figure 8 also shows the predicted spectra of the errors due to system thermal noise and the ionosphere. The measured spectrum is in good agreement with the predicted system noise effect at the highest frequencies. The ionospheric error spectrum was computed assuming white frequency noise for lower frequencies and Kolmogorov turbulence for higher frequencies, with a magnitude consistent with the error budget. The ionosphere may be affecting the measured spectrum at the lowest frequencies.

The measured spectrum of SBI data residuals and the spectra of predicted error sources are shown in Fig. 9 for

its medium gain antenna, and the predicted solar plasma and ionospheric errors are less for February 20, since the angular separation between the two spacecraft was smaller during the February 20 SBI data arc. The measured spectrum is again in qualitative agreement with the prediction, which suggests that error models have been successfully applied for two different observation geometries.

The SBI measurement error as a function of averaging interval is displayed in Fig. 10. Curves are drawn for the predicted error for both the August 1990 and the February 1991 data sets. Points are plotted for the errors as measured for the August 11 and the February 20 data arcs. The measured values are about 25 percent less than predicted for 5-min averages, and about a factor of three less than predicted for shorter averaging intervals. Good characterization of measurement errors is more important at the longer averaging intervals, since orbit determination will be affected by the errors which remain after data averaging.

It should be noted that in acquiring the SBI data, the signals are not simultaneously recorded; the spacecraft signals from the two spacecraft are recorded in separate channels which are multiplexed with 1-sec dwells. Time multiplexing of the channels can conceivably introduce additional instrumental phase errors that do not cancel when the spacecraft signals are doubly differenced. For this reason, several phase calibration tones are injected near the beginning of the instrumental path into each channel. The phase calibration tone generator is locked to the station's hydrogen maser, providing frequency stability at the level of 10^{-15} sec/sec for 1000-sec intervals. Calculation of the power spectrum of doubly differenced phase calibration tone residuals for the August 1990 data set showed that the errors contributed by imperfect instrumental cancellation between channels was insignificant; the power spectrum of doubly differenced phase calibration tone residuals was

several orders of magnitude below the measured spectrum illustrated in Fig. 8.

VI. Discussion

The SBI data acquired in 1990-1991 from the Magellan and Pioneer 12 orbiters at Venus have provided an initial assessment of SBI measurement errors. The scatter in the residuals is consistent with predictions for averaging intervals of 5 min, and less than predicted by about a factor of 3 for shorter averaging intervals. The error budget for these measurements, made at S-band, is entirely dominated by solar plasma fluctuations at the shorter averaging intervals. The plasma error model is based on line-of-sight observations. The eventual availability of dual-frequency same-beam radio metric data is expected to provide substantial improvement of this model for predicting SBI errors.

A single S-band SBI measurement (5-min average) determines spacecraft plane-of-sky separation with an accuracy of 40 to 100 m. Measurements made at X-band are expected to be more accurate by an order of magnitude. Measurement errors other than solar plasma are not distinguishable in the available S-band data set. Additional data acquisition is anticipated to further examine measurement system errors.

Demonstration of improved orbital accuracy is essential to validate the performance of the SBI measurement technique. Orbits generated using disjointed data arcs, and various combinations of Doppler, differenced Doppler, and SBI data, may be propagated to a common epoch and compared. The radio metric data in hand from Magellan and Pioneer 12 provide the opportunity to demonstrate the contribution of the SBI technique toward orbital accuracy improvement.

Acknowledgments

The authors thank the Magellan and Pioneer projects for their cooperation and support in obtaining the radio metric data analyzed here. The reference spacecraft trajectories used in this analysis were supplied by the Magellan and Pioneer navigation teams.

References

- [1] J. R. Smith and R. Ramos, "Data Acquisition for Measuring the Wind on Venus from Pioneer Venus," *IEEE Transactions on Geoscience and Remote Sensing*, vol. GE-18, no. 1, pp. 126-130, January 1980.
- [2] R. A. Preston, C. E. Hildebrand, G. H. Purcell, Jr., J. Ellis, C. T. Stelzried, S. G. Finley, R. Z. Sagdeev, V. M. Linkin, V. V. Kerzhanovich, V. I. Altunin, L. R. Kogan, V. I. Kostenko, L. I. Matveenko, S. V. Pogrebenko, I. A. Strukov, E. L. Akim, Yu. N. Alexandrov, N. A. Armand, R. N. Bakitko, A. S. Vyshlov, A. F. Bogomolov, Yu. N. Gorchankov, A. S. Selivanov, N. M. Ivanov, V. F. Tichonov, J. E. Blamont, L. Boloh, G. Laurans, A. Boisshot, F. Biraud, A. Ortega-Molina, C. Rosolen, and G. Petit, "Determination of Venus Winds by Ground-Based Radio Tracking of the VEGA Balloons," *Science*, vol. 231, pp. 1414-1416, March 21, 1986.
- [3] C. C. Counselman III, H. F. Hinteregger, and I. I. Shapiro, "Astronomical Applications of Differential Interferometry," *Science*, vol. 178, pp. 607-608, November 10, 1972.
- [4] R. W. King, C. C. Counselman III, and I. I. Shapiro, "Lunar Dynamics and Selenodesy: Results from Analysis of VLBI and Laser Data," *J. Geophys. Res.*, vol. 81, no. 35, pp. 6251-6256, December 10, 1976.
- [5] J. S. Border and R. D. Kahn, "Relative Tracking of Multiple Spacecraft by Interferometry," in *Advances in the Astronautical Sciences: Orbital Mechanics and Mission Design*, vol. 69, edited by J. Teles, San Diego, California: Univelt, 1989.
- [6] J. S. Border and W. M. Folkner, "Differential Spacecraft Tracking by Interferometry," *Proceedings of the CNES International Symposium on Space Dynamics*, Toulouse, France, November 1989.
- [7] W. M. Folkner and J. S. Border, "Orbiter-Orbiter and Orbiter-Lander Tracking Using Same-Beam Interferometry," *TDA Progress Report 42-109*, vol. January-March 1992, Jet Propulsion Laboratory, Pasadena, California, pp. 74-86, May 15, 1992.
- [8] R. D. Kahn, W. M. Folkner, C. D. Edwards, and A. Vijayaraghavan, "Position Determination of a Lander and Rover at Mars With Earth-Based Differential Tracking," *TDA Progress Report 42-108*, vol. October-December 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 279-293, February 15, 1992.
- [9] D. B. Engelhardt, J. B. McNamee, S. K. Wong, F. G. Bonneau, E. J. Graat, R. J. Haw, G. R. Kronschnabl, and M. S. Ryne, "Determination and Prediction of Magellan's Orbit," paper AAS-91-180, AAS/AIAA Spaceflight Mechanics Meeting, Houston, Texas, February 11-13, 1991.
- [10] W. G. Melbourne and D. W. Curkendall, "Radio Metric Direction Finding: A New Approach to Deep Space Navigation," paper presented at AAS/AIAA Astrodynamics Specialist Conference, Jackson Hole, Wyoming, September 7-9, 1977.
- [11] S. W. Thurman, "Deep-Space Navigation with Differenced Data Types, Part II: Differenced Doppler Information Content," *TDA Progress Report 42-103*, vol. July-September 1990, Jet Propulsion Laboratory, Pasadena, California, pp. 61-69, November 15, 1990.

- [12] K. M. Liewer, "DSN Very Long Baseline Interferometry System Mark IV-88," *TDA Progress Report 42-99*, vol. January-March 1988, Jet Propulsion Laboratory, Pasadena, California, pp. 239-246, May 15, 1988.
- [13] J. S. Border and E. R. Kursinski, "Deep Space Tracking and Frequency Standards," *Proceedings of the 45th Annual Symposium on Frequency Control, IEEE* Catalog No. 91CH2965-2, Los Angeles, California, May 29-31, 1991.
- [14] C. Dunn, S. Lichten, D. Jefferson and J. S. Border, "Sub-Nanosecond Clock Synchronization and Precision Deep Space Tracking," *Proceedings of the Twenty-Third Annual Precise Time and Time Interval (PTTI) Applications and Planning Meeting*, Pasadena, California, December 3-5, 1991.
- [15] N. C. Ham, "VLBI System (BLK I) IF-Video Down Conversion Design," *TDA Progress Report 42-79*, vol. July-September 1984, Jet Propulsion Laboratory, Pasadena, California, pp. 172-188, November 15, 1984.
- [16] A. P. Freedman, "Combining GPS and VLBI Earth-Rotation Data for Improved Universal Time," *TDA Progress Report 42-105*, vol. January-March 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 1-12, May 15, 1991.
- [17] U. J. Lindqwister, A. P. Freedman, and G. Blewitt, "A Demonstration of Centimeter-Level Monitoring of Polar Motion With the Global Positioning System," *TDA Progress Report 42-108*, vol. October-December 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 1-9, February 15, 1992.
- [18] M. H. Finger and W. M. Folkner, "A Determination of the Radio-Planetary Frame Tie From Comparison of Earth Orientation Parameters," *TDA Progress Report 42-109*, vol. January-March 1992, Jet Propulsion Laboratory, Pasadena, California, pp. 1-21, May 15, 1992.
- [19] R. N. Treuhaft and G. E. Lanyi, "The Effect of the Dynamic Wet Troposphere on Radio Interferometric Measurements," *Radio Science*, vol. 22, no. 2, pp. 251-265, March-April 1987.
- [20] R. Woo and J. W. Armstrong, "Spacecraft Radio Scattering Observations of the Power Spectrum of Electron Density Fluctuations in the Solar Wind," *J. Geophys. Res.*, vol. 84, no. A12, pp. 7288-7296, December 1, 1979.
- [21] R. D. Kahn and J. S. Border, "Precise Interferometric Tracking of Spacecraft at Low Sun-Earth-Probe Angles," paper AIAA-88-0572, AIAA 26th Aerospace Sciences Meeting, Reno, Nevada, January 11-14, 1988.
- [22] R. Woo, F.-C. Yang, K. W. Yip, and W. B. Kendall, "Measurements of Large-Scale Density Fluctuations in the Solar Wind Using Dual-Frequency Phase Scintillations," *Astrophysical J.*, vol. 210, pp. 568-574, December 1, 1976.

Table 1. Pioneer 12 and Magellan transmitter frequencies.

Transmitter frequency, MHz	S-band	X-band
Radio frequency	2293.81	8410.63
Frequency offset	4.16	0.11

Table 2. SBI measurement error model assumptions.

Source	Magnitude
Station clock offset	0.1 μ sec
Zenith troposphere error	4 cm
Zenith ionosphere error	5 TEC units
Baseline component error*	21.2 cm

* Includes station location, Earth orientation, and radio-planetary frame tie.

Table 3. SBI observation geometry.

Component	August 1990	February 1991
Venus		
Right ascension, deg	120.3	355.7
Declination, deg	20.8	-3.2
Sun-Earth-Venus angle, deg	21.3	25.8
Distance from Earth, AU	1.563	1.472
Elevation angle, deg		
Madrid	15.	—
Goldstone	45.	20.
Canberra	15.	35.
Differential elevation angle, mrad		
Madrid	0.16	—
Goldstone	0.30	0.07
Canberra	0.16	0.07
Angular separation, mrad		
Goldstone-Madrid	0.30	—
Goldstone-Canberra	0.24	0.09
Change in S-band Doppler shift over data arc, kHz		
Magellan	47.	54.
Pioneer 12	2.4	67.

Table 4. PVO and MGN link analysis, for PVO S-band and X-band carriers, MGN S-band carrier from medium-gain antenna (MGA) and high-gain antenna (HGA), and MGN -16th harmonic of high-rate telemetry subcarrier at X-band.

Link component	PVO S-band	PVO X-band	MGN S-band MGA	MGN S-band HGA	MGN X-band
Power transmitted, dBm	40.0	27.8	35.5	35.6	42.4
Spacecraft antenna gain, dB	25.2	35.0	18.7	35.9	47.9
Space loss, 1.5 AU, dB	-267.1	-278.3	-267.1	-267.1	-278.3
P_{Tone}/P_{Total} , dB	-8.4	0.0	-3.5	-7.4	-38.1
Polarization loss, dB	0.0	0.0	-0.5	-3.0	0.0
Receiving antenna gain, dB	56.0	68.0	56.0	56.0	68.0
P_{Tone} , dBm	-154.3	-147.5	-160.9	-150.0	-158.1
N_0^a , dBm/Hz	-182.1	-184.6	-182.1	-182.1	-184.6
P_{Tone}/N_0 , dB·Hz	27.8	37.1	21.2	32.1	26.5
P_{Tone}/N_0 , W/W	603.	5130.	132.	1620.	447.
$SNR_{V,1sec}$	27.7	80.8	13.0	45.4	23.9

^a System temperature 45 K at S-band and 25 K at X-band.

Table 5. SBI data acquired from Magellan and Pioneer 12.

Session	Goldstone-Madrid, hr	Goldstone-Canberra, hr	Radio frequency, GHz
August 1990	2	3	2.3
February 1991	0	8	2.3
April 1991	7	13.5	2.3

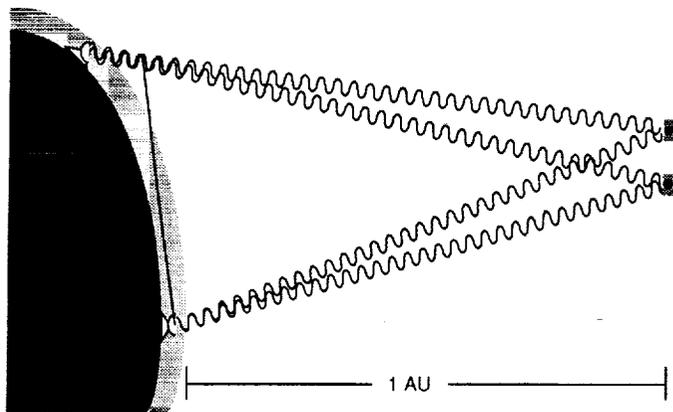


Fig. 1. Geometry of SBI measurements.

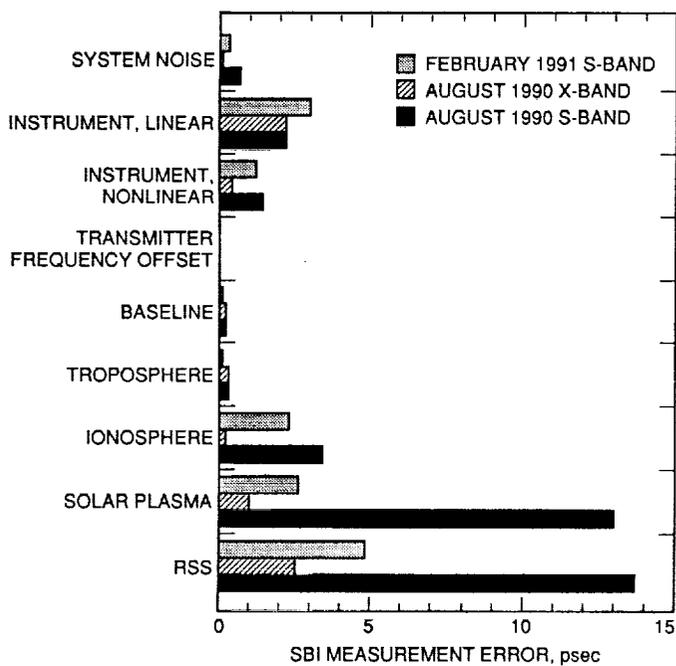


Fig. 2. Predicted SBI measurement errors for MGN-PVO for 5-min averages. Predicted X-band errors are shown for comparison only.

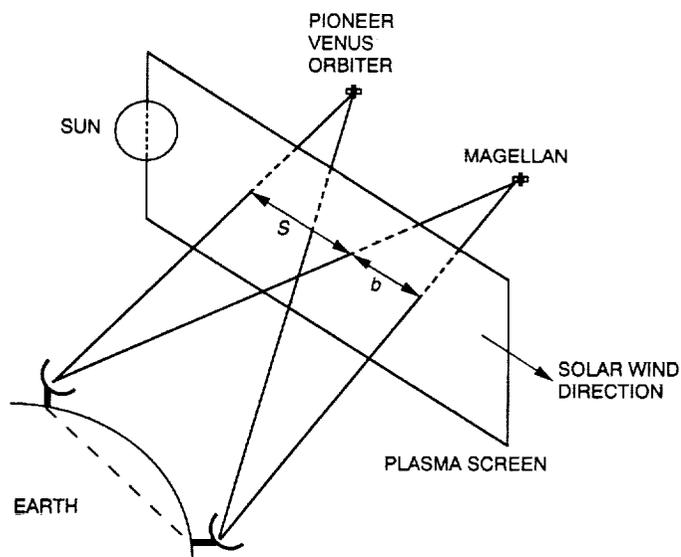


Fig. 3. Geometry of radio signals traversing the solar plasma.

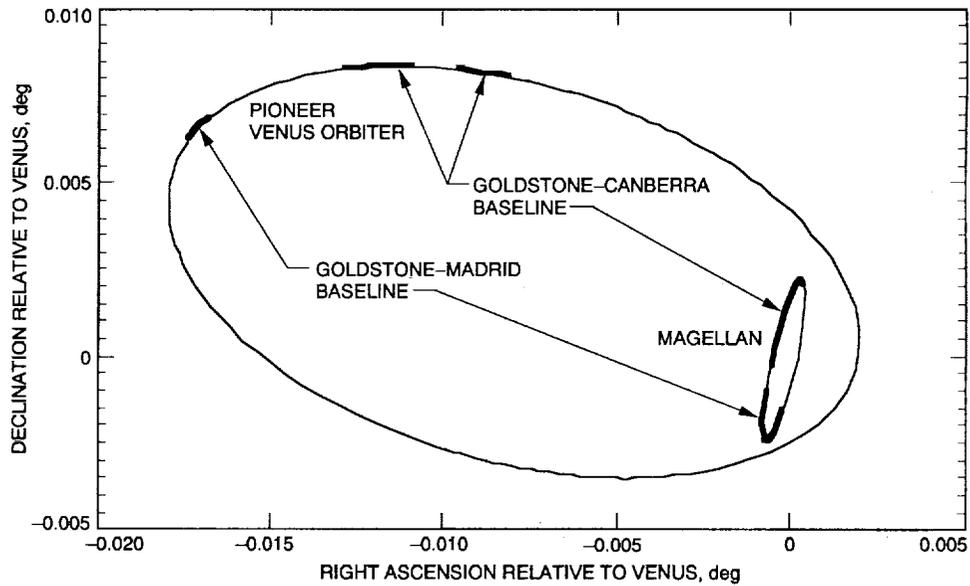


Fig. 4. Venus-centered orbit traces of MGN and PVO as seen from Earth, for geometry of August 1990.

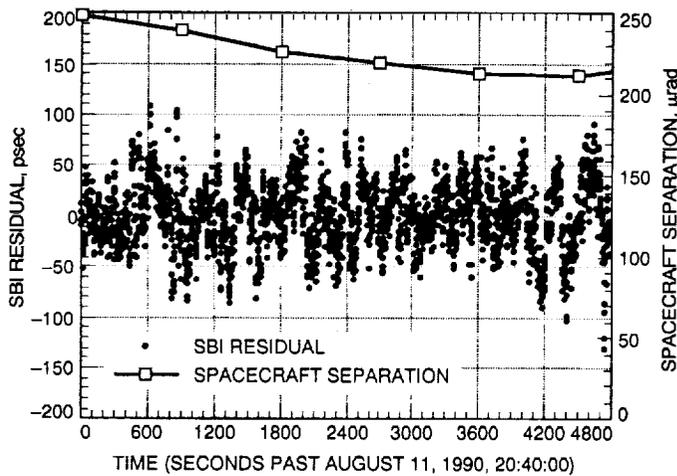


Fig. 5. Two-second MGN-PVO SBI residuals (S-band) for Goldstone-Canberra, August 11, 1990.

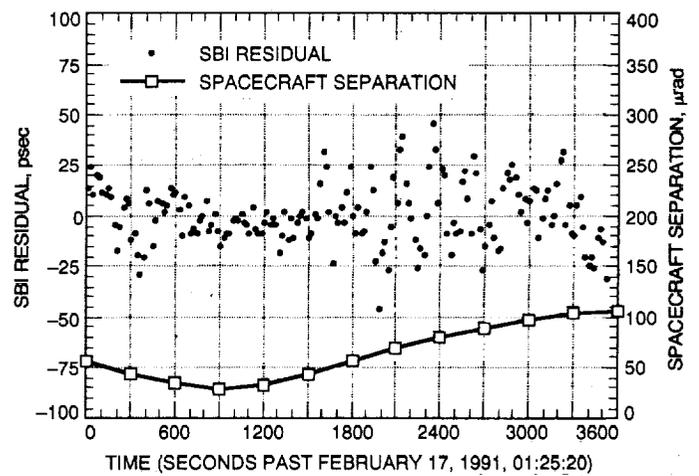


Fig. 6. Twenty-second MGN-PVO SBI residuals (S-band) for Goldstone-Canberra, February 17, 1991.

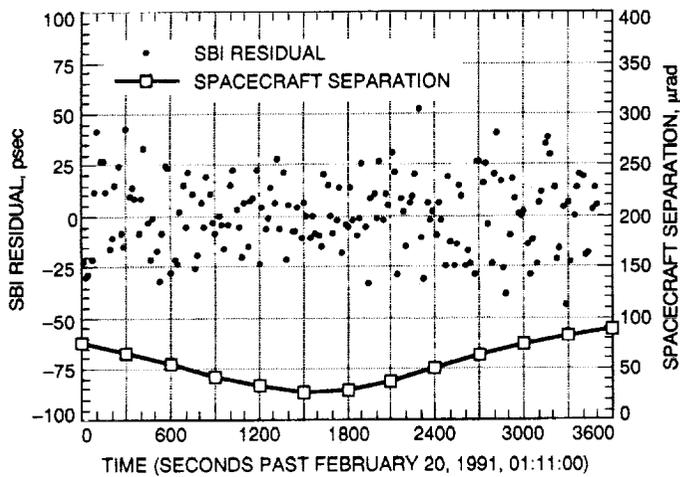


Fig. 7. Twenty-second MGN-PVO SBI residuals (S-band) for Goldstone-Canberra, February 20, 1991.

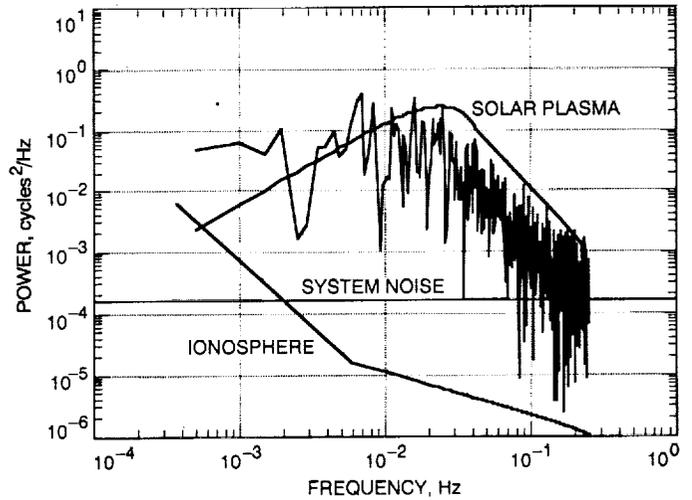


Fig. 9. Power spectrum of doubly differenced residual phase (S-band), observed and predicted, for Goldstone-Canberra, February 20, 1991.

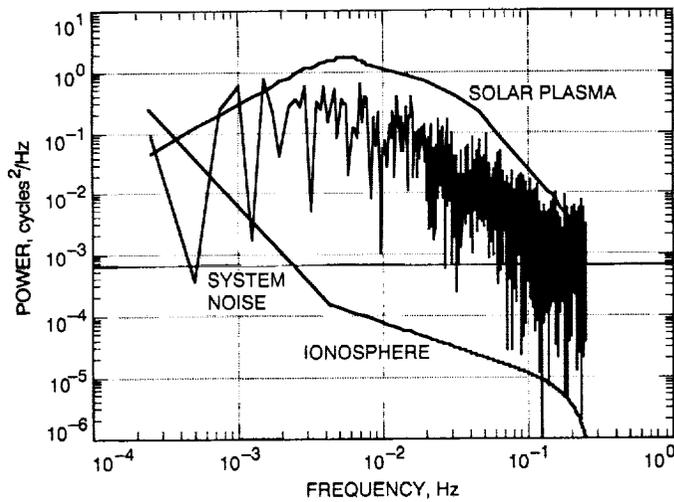


Fig. 8. Power spectrum of doubly differenced residual phase (S-band), observed and predicted, for Goldstone-Canberra, August 11, 1990.

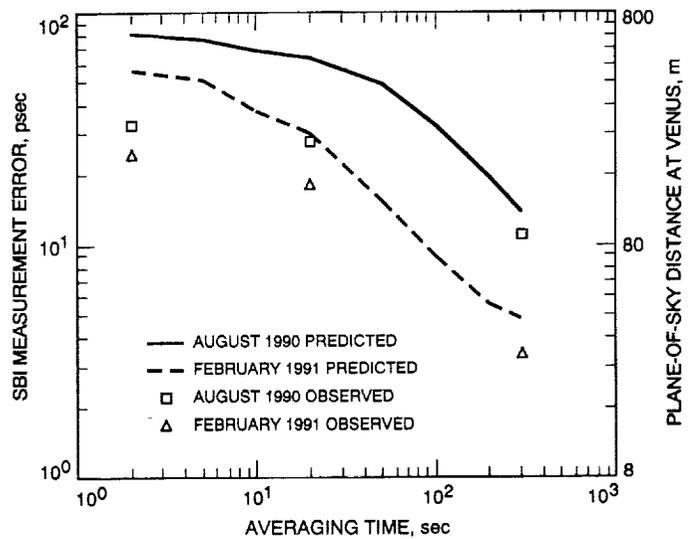


Fig. 10. Observed and predicted SBI measurement error (S-band) as a function of averaging interval.

52-13
128435
N93-19415

Application of High-Precision Two-Way Ranging to Galileo Earth-1 Encounter Navigation

V. M. Pollmeier and S. W. Thurman
Navigation Systems Section

The application of precision two-way ranging to orbit determination with relatively short data arcs is investigated for the Galileo spacecraft's approach to its first Earth encounter (December 8, 1990). Analysis of previous S-band (2.3-GHz) ranging data acquired from Galileo indicated that under good signal conditions sub-meter precision and 10-m ranging accuracy were achieved. It is shown that ranging data of sufficient accuracy, when acquired from multiple stations, can sense the geocentric angular position of a distant spacecraft. A range data filtering technique, in which explicit modeling of range measurement bias parameters for each station pass is utilized, is shown to largely remove systematic ground system calibration errors and transmission media effects from the Galileo range measurements, which would otherwise corrupt the angle-finding capabilities of the data. The accuracy of the Galileo orbit solutions obtained with S-band Doppler and precision ranging were found to be consistent with simple theoretical calculations, which predicted that angular accuracies of 0.26–0.34 μ rad were achievable. In addition, the navigation accuracy achieved with precision ranging was marginally better than that obtained using delta-differenced one-way range (Δ DOR), the principal data type that was previously used to obtain spacecraft angular position measurements operationally.

I. Introduction

The approach phase leading up to the Galileo spacecraft's first Earth encounter (designated Earth-1) provided a good opportunity to test the viability of high-precision two-way ranging as an operational radio metric data type. Two-way ranging data acquired by Deep Space Network (DSN) stations have been accurate to 15 m or better for nearly 20 years, depending upon the frequency band and station-spacecraft radio link characteristics. Such data have typically been utilized for orbit determination at assumed accuracies of 100–1000 m, due to the effects of station delay and transmission media calibration errors,

and the influence of small, poorly modeled spacecraft non-gravitational forces. Since the early 1970s, evolutionary improvements in the accuracy and stability of timing systems, station delay calibration procedures, and transmission media calibration techniques, coupled with more sophisticated orbit determination software, now make it possible to reconsider the use of precision ranging for interplanetary spacecraft navigation.

In a recent experiment conducted with radio metric data acquired from the Ulysses spacecraft, two-way ranging data were processed with a new range data-filtering

technique that made it possible to successfully utilize the data at an assumed accuracy of 10 m for the first time [1]. This filtering technique utilized estimated parameters to explicitly account for and remove residual ground system calibration errors and solar plasma-induced delays from the ranging data. Other factors that contributed to the success of the Ulysses experiment were the accuracy and consistency of the DSN station delay calibrations and the utilization of a new DSN station location set developed by Folkner and Dewey.¹

The use of Galileo ranging for a second test of the range data-filtering technique used to process the Ulysses data was motivated by the earlier results of the Galileo Venus encounter orbit solutions. During Galileo's approach to Venus and its subsequent flyby, good signal strength was obtained from the spacecraft's two low-gain S-band (2.3-GHz) antennas, yielding point-to-point two-way range noise (indicative of the precision of the data) of under 1 m and an apparent accuracy of 10 m or better [2]. In addition, the station delay calibrations during this time period appeared to be of good consistency, showing little variation over the month prior to the encounter. This article describes an investigation which reexamined the orbit determination that was utilized for the design of the last targeting maneuver prior to the Earth-1 encounter. A brief discussion of the theoretical basis for the ability of high-accuracy ranging data to sense spacecraft angular position is presented, as well as comparisons of different solutions obtained using combinations of various data types, including two-way Doppler and ranging, and Δ DOR.

II. Theoretical Background

A simple investigation of the ability of ranging and Doppler data to determine the trajectory of a distant spacecraft can be conducted by analyzing the theoretical precision with which the geocentric spacecraft motion can be sensed from one or two passes of data. Similar analyses have been performed previously for range and Doppler data separately [1,3,4]. The station-spacecraft tracking geometry is illustrated in Fig. 1. The topocentric range, ρ , and range-rate, $\dot{\rho}$, can be accurately approximated over short periods of time (up to roughly 24 hr) in terms of the geocentric spacecraft range (r), range rate (\dot{r}), declination (δ), and right ascension (α), as follows:

$$\rho \approx r - (r_s \cos \delta \cos H + z_s \sin \delta) \quad (1)$$

$$\dot{\rho} \approx \dot{r} + \omega r_s \cos \delta \sin H \quad (2)$$

where

r_s = station distance from Earth's spin axis (spin radius)

z_s = station height above Earth's equator (z -height)

ω = Earth rotation rate (7.3×10^{-5} rad/sec)

$H = \alpha_g + \lambda - \alpha$

and

α_g = right ascension of Greenwich meridian

λ = station east longitude

From Eqs. (1) and (2), it can be seen that four of the six components of the geocentric spacecraft trajectory (r , \dot{r} , δ , and α) can be sensed by range and range-rate measurements. Over the time period of interest, \dot{r} , δ , and α are nearly constant. A determination of the remaining two coordinates, δ and α , and hence the complete trajectory, normally requires the acquisition of multiple passes of data over a period of several days. The accumulated information in each ranging and Doppler pass can be thought of as a multidimensional measurement of the spacecraft trajectory, with the statistical combination of several such measurements yielding a complete determination of the trajectory.

A simple least-squares error analysis of estimates of r , \dot{r} , δ , and α , derived from a single pass of range and Doppler data, can be formulated analytically (refer to the paper by Hamilton and Melbourne [3] for more details). For the purposes of this analysis, it is assumed that \dot{r} , δ , and α are constants, and that r varies linearly with time. The information matrix, J , for these coordinates, assuming a tracking pass in which the station-spacecraft hour angle H varies as $-\psi \leq H \leq +\psi$, can be expressed as

$$\begin{aligned} J \approx & \left(\frac{1}{\sigma_\rho^2 \omega \Delta t} \right) \int_{-\psi}^{+\psi} [\partial \rho / \partial (r, \dot{r}, \delta, \alpha)]^T \\ & \times [\partial \rho / \partial (r, \dot{r}, \delta, \alpha)] dH \\ & + \left(\frac{1}{\sigma_{\dot{\rho}}^2 \omega \Delta t} \right) \int_{-\psi}^{+\psi} [\partial \dot{\rho} / \partial (r, \dot{r}, \delta, \alpha)]^T \\ & \times [\partial \dot{\rho} / \partial (r, \dot{r}, \delta, \alpha)] dH \end{aligned} \quad (3)$$

¹ W. M. Folkner and R. J. Dewey, "Radio Source Catalog and Station Location Set for Ulysses," JPL Interoffice Memorandum 335.1-90-048, Jet Propulsion Laboratory, Pasadena, California, September 13, 1990.

where

σ_ρ = range measurement noise one-sigma uncertainty

$\sigma_{\dot{\rho}}$ = range-rate (Doppler) measurement noise one-sigma uncertainty

Δt = time interval between measurements

In Eq. (3), it is assumed that the time between measurements, Δt , is the same for both the range and Doppler measurements. The partial derivatives appearing in Eq. (3) at some time t , with respect to the geocentric coordinates at time t_0 , where t_0 is assumed to be the time at which the spacecraft crosses the local meridian of the tracking station, are as follows:

$$\frac{\partial \rho}{\partial(r, \dot{r}, \delta, \alpha)} \approx [1, t - t_0, r_s \sin \delta \cos H - z_s \cos \delta, -r_s \cos \delta \sin H] \quad (4)$$

$$\frac{\partial \dot{\rho}}{\partial(r, \dot{r}, \delta, \alpha)} \approx [0, 1, -\omega r_s \sin \delta \sin H, -\omega r_s \cos \delta \cos H] \quad (5)$$

The error covariance, Λ , for r , \dot{r} , δ , and α at time t_0 is simply

$$\Lambda = J^{-1} = \begin{bmatrix} \sigma_r^2 & 0 & \sigma_{r\delta}^2 & 0 \\ 0 & \sigma_{\dot{r}}^2 & 0 & \sigma_{\dot{r}\alpha}^2 \\ \sigma_{r\delta}^2 & 0 & \sigma_\delta^2 & 0 \\ 0 & \sigma_{\dot{r}\alpha}^2 & 0 & \sigma_\alpha^2 \end{bmatrix} \quad (6)$$

where

$$\sigma_\delta^2 = \frac{\omega \Delta t}{(r_s \sin \delta)^2} f_1(\psi, \sigma_\rho^2, \sigma_{\dot{\rho}}^2) \quad (7)$$

$$\sigma_\alpha^2 = \frac{\omega \Delta t}{(r_s \cos \delta)^2} f_2(\psi, \sigma_\rho^2, \sigma_{\dot{\rho}}^2) \quad (8)$$

Equations (7) and (8) are similar to expressions derived by Anderson [4] in an earlier analysis of this same problem (the functions f_1 and f_2 are not shown explicitly, as they are rather complex). As noted by Anderson, σ_δ is proportional to $1/\sin \delta$, and will theoretically become infinite for spacecraft located on the celestial equator ($\delta = 0$). Hamilton and Melbourne [3] found an equivalent result for a single pass of Doppler data only. In contrast, σ_α is seen

from Eq. (8) to be proportional to $1/\cos \delta$, which has little (± 10 percent) variation over the declination range spanned by the ecliptic plane (± 24 deg), in which most interplanetary spacecraft trajectories lie. Although formulas for σ_r and $\sigma_{\dot{r}}$ are not explicitly given, Eq. (6) predicts that these quantities are determined with a precision limited only by the precision of the range and Doppler measurements and the number of measurements acquired. Thus, for the case of single-station tracking, the ability of ranging and Doppler data to determine the spacecraft declination depends heavily on the tracking geometry.

The situation described above changes dramatically when an additional pass of ranging and Doppler data from a properly chosen second station is added into the information matrix. Consider a scenario in which a tracking pass is acquired from a station with z -height, z_s , and spin radius, r_s , followed immediately by another pass from a second station, with z -height, $-z_s$, and spin radius, r_s . This choice of station coordinates is not arbitrary; stations located at the DSN complexes at Goldstone and near Canberra have spin radii that are nearly equal (to within about 5 km) and z -heights that are nearly equal in magnitude but have opposite signs. Applying the assumptions used in the single-pass analysis to this case yields an error covariance matrix, Λ , that incorporates the information matrix obtained from the first pass, designated J_1 , and the information matrix from the second pass, J_2 :

$$\Lambda = [J_1 + J_2]^{-1} \quad (9)$$

For the case of $\delta = 0$, the formula for σ_δ obtained from Eq. (9) reduces to a simple form

$$\sigma_\delta = \frac{\sigma_\rho}{2z_s} \sqrt{\frac{\omega \Delta t}{\psi}} \quad (10)$$

From Eq. (10), it can be seen that the z -height component of the baseline formed by the two stations enables a determination of δ , and that this determination is provided solely by the ranging data. The result for σ_α obtained in Eq. (10) is simply equal to the expression for σ_α from Eq. (8) multiplied by a factor of $1/\sqrt{2}$.

A simplified illustration of the result obtained in Eq. (10) is shown in Fig. 2, for a spacecraft at near-zero declination (i.e., $\sin \delta \approx \delta$) being tracked by two stations located on a two-dimensional Earth. In Fig. 2, the two stations shown have z -heights equal in magnitude but opposite in sign, as was assumed in developing Eqs. (9) and

(10). The spacecraft declination is sensed through the difference between the range measured from the two different stations. As explained by Taylor et al. [5], the greatest accuracy in determining this range difference is achieved by explicitly differencing simultaneous or near-simultaneous range measurements obtained during periods of mutual visibility. If the spacecraft dynamics and the range measurements can be modeled with sufficient accuracy, though, this explicit differencing is not required; the two-way ranging data from the different station passes implicitly contain the information needed to determine δ , as shown in Eq. (10) above.

Using Eq. (10), the angular precision that can theoretically be achieved by using two passes of S-band (2.3-GHz) ranging and Doppler data for $\delta = 0$ was computed and plotted in Fig. 3, as a function of the combined tracking time from the two stations, which were assumed to have r_s and z_s coordinates corresponding to the DSN Goldstone and Canberra complexes. The assumed ranging and Doppler measurement accuracies used to construct Fig. 3 ($\sigma_\rho = 10$ m and $\sigma_\delta = 1$ mm/sec) are based on previous experience with Galileo S-band ranging and Doppler data [2]. In Fig. 3, the total tracking time is assumed to be divided equally between the two participating stations. During the Galileo spacecraft's approach to the Earth-1 encounter, typical tracking pass lengths were 9 to 10 hr; Fig. 3 indicates that for two 10-hr passes, the theoretical angular precision achieved is about $0.26 \mu\text{rad}$ in declination and about $0.34 \mu\text{rad}$ in right ascension. Subsequent calculations using Eq. (10) for nonzero values of δ ranging from -24 to $+24$ deg (not shown here) yielded results that were within 10 percent of the data shown in Fig. 3. In comparison, the angular precision of Galileo S-band ΔDOR , the principal data type used to obtain angular measurements during actual Earth-1 encounter navigation operations, was about 0.04 – $0.08 \mu\text{rad}$, depending upon the tracking geometry [6].

Although not as accurate as ΔDOR in this experiment, two-way ranging is a much simpler data type to employ operationally, in that the data are easier to acquire and process. In addition, it will be shown subsequently that the superior angular precision of ΔDOR does not always translate into an equivalent level of navigation accuracy. In a typical mission operations environment, the scheduling and postprocessing requirements associated with ΔDOR result in data acquisitions every 1 to 2 days at best, and more often at intervals of 3 to 7 days instead (e.g., 19 ΔDOR measurements were acquired over the 44-day data arc used in this experiment; 2750 two-way range measurements were acquired during the same period). This sparsity of ΔDOR data sometimes leads to navigation ac-

curacies that are, in an angular sense, poorer than the theoretical angle-finding capability of the data.

While the theoretical results above show that ranging can overcome the dependence of Doppler-based angle determination on the tracking geometry, it must be recognized that the effects of systematic range measurement errors, principally station delay calibration errors, will not necessarily be reduced through statistical averaging, as will the effects of random errors. These systematic errors must be accounted for in some way, or reduced a priori through the use of very accurate calibrations.

III. Station Delay Calibrations

To account for the effects of systematic bias errors on the ranging data, pass-specific bias parameters were estimated for each Deep Space Station (DSS) used to acquire ranging from Galileo. In addition to station delay calibration errors, which generally do not change very much over the duration of a single pass, these bias parameters were intended to represent slowly varying range delays due to the solar plasma, which are often the largest nongeometric component of range measurements acquired with S-band uplink and downlink frequencies [1]. The data were divided into batches so that no two ranging passes from the same station were included in a single batch. Stochastic range bias parameters were then estimated for each station during each successive batch by using a batch-sequential filter algorithm.

An examination of the station delay calibration data obtained during the Earth-1 approach phase indicated that for the 60 days prior to the encounter, the values of station delay calibrations for the DSN 70-meter subnet (most of the Galileo Earth-1 approach radio metric data were acquired from this subnet) were very consistent and showed little day-to-day variation. Great effort was made on the part of DSN station personnel to maintain the consistency of the 70-m station configurations during this phase of the mission. The standard deviation of the station delay calibrations was observed to be just 55 cm, with the largest variations being on the order of 1.5 m. Since tracking passes were obtained infrequently from the 34-m standard (STD) subnet, the sparsity of station delay calibration data from this subnet prevented any similar analysis. The Sun-Earth-spacecraft angle was quite large within the data arc (greater than 150 deg); therefore, the anticipated magnitude of solar plasma-induced delays in the S-band range measurements was 1 m or less, assuming an average level of solar activity [7]. Based on these considerations, the stochastic range biases associated with the 70-m stations were assigned a priori uncertainties of 2 m, and the

range biases associated with the 34-m STD stations were assigned a priori uncertainties of 10 m.

IV. Analysis

The Earth-1 orbit determination analysis for this experiment consisted of recomputing the orbit determination delivery that was used for the design of the final Earth-targeting maneuver and used several different data sets and assumed data accuracies. The data arc used for the solutions extended from October 10, 1990 (59 days prior to encounter) to November 23, 1990 (15 days prior to encounter). This time period corresponds to Earth-spacecraft distances ranging from 50–12.5 million km, and a geocentric spacecraft declination of 15–13 deg. The radio metric data acquired included 3740 Doppler points (600-sec count time) and 2750 range points. Additionally, 19 Δ DOR observations were obtained, including 11 observations from the DSN Goldstone–Canberra baseline, and 8 observations from the Goldstone–Madrid baseline. Table 1 summarizes the parameters and assumptions that were used in the orbit determination filter model. In Table 1, the estimated parameters are those that were explicitly solved for in the estimation process; the consider parameters were not estimated, but the effects of uncertainty in these quantities was accounted for (i.e., “considered”) when calculating the error covariance associated with the solution for the estimated parameters. Also in Table 1, the radial and transverse components of the solar radiation pressure model refer to the direction parallel to the Sun–spacecraft line, and the two directions orthogonal to that line, respectively.

For the set of solutions that was computed, several different choices of data set and data weighting (i.e., specification of the assumed measurement noise level for each data type) were exercised in order to determine the effect of each variation on the predicted aim point for the encounter. These solutions were then compared with a highly accurate (50-m) post-flyby reconstruction of the trajectory that was computed using both pre- and post-encounter radio metric data. For the precision ranging analysis, a range data weight of 10 meters was used. Although the noise level previously observed in Galileo ranging data was at the submeter level, a weight of 10 m was chosen in light of the presence of 1- to 2-meter-level systematic ionospheric calibration errors that could affect data acquired at S-band frequencies. For comparison purposes, a range data weight of 1 km was used in several solutions, as this value is representative of more traditional methods of utilizing ranging (1 km was, in fact, the range weight used operationally for the Earth-1 encounter). Two sets of solutions were

calculated; in the first set a Doppler weight of 1 mm/sec (60-sec count time) was used for all solutions, and in the second set a Doppler weight of 2 mm/sec (60-sec count time) was employed. The Doppler data weight used during actual Earth-1 encounter operations was 1 mm/sec, which is commensurate with the inherent accuracy of the data. In each set, solutions were constructed using Doppler data only; Doppler plus 1-km range; Doppler, 1-km range and 50-cm Δ DOR; and Doppler plus 10-m range. A final solution was constructed using only 10-m range for comparison purposes. The stochastic range bias filter model was employed in all cases involving 10-m range.

V. Results

The results of the analysis are shown in Figs. 4 and 5, and are summarized in Table 2. Figures 4 and 5 portray the two sets of solutions obtained with Doppler weights of 1 mm/sec and 2 mm/sec, respectively, in an Earth-centered aiming plane coordinate system.² The one-sigma dispersion ellipses associated with each solution (representing a 39-percent confidence interval) are also shown in Figs. 4 and 5. For the Earth-1 encounter, the aiming plane was nearly coincident with the plane of the sky, that plane which is normal to the Earth–spacecraft line-of-sight, over the entire data arc. Therefore, the ability of different radio metric data types to determine the aim point for the encounter was very closely related to their ability to measure the geocentric spacecraft angular position over the time span of the data arc. Thus, this encounter represents a fairly direct test of the angle-finding capability of S-band precision ranging data.

The solution utilizing 1-mm/sec Doppler and high-precision ranging, shown in Fig. 4, resulted in the closest agreement with the post-flyby reconstruction of all the solutions performed, with an error of 4.2 km in the aiming plane. This error translates into an angular error of 0.34 μ rad, a value that is in good agreement with the theoretical precision of two Northern and Southern Hemisphere 9- or 10-hr DSN tracking passes, shown in Fig. 3. This precision range result compares quite favorably with the 1-mm/sec Doppler, 1-km range, 50-cm Δ DOR solution (also shown in Fig. 4), which had an error of 7.8 km

² The aiming plane or “B-plane” coordinates system is defined by three unit vectors, **S**, **T**, and **R**; **S** is parallel to the incoming asymptotic velocity vector, **T** is parallel to the ecliptic plane (mean plane of the Earth’s orbit), and **R** completes an orthogonal triad with **S** and **T**. The aim point for a planetary flyby is defined by the miss vector, **B**, which lies in the **T**–**R** plane, and can be thought of as specifying where the point of closest approach would be if the target planet had no mass and did not deflect the flight path.

(equivalent to $0.62 \mu\text{rad}$); this solution was the best solution obtained during actual Earth-1 encounter operations. The relatively poor performance of ΔDOR is attributed to the number of velocity changes (due to attitude update maneuvers and propellant line flushings) that had to be estimated by the orbit determination filter, and the sparsity of the ΔDOR data set (19 ΔDOR points versus 2750 range points). Operational complexities associated with the scheduling of ΔDOR and the sequencing procedures used by Galileo made it difficult to acquire a large data set for the Earth-1 encounter. In addition, ΔDOR scheduling difficulties during the Earth-1 approach resulted in a somewhat irregular distribution of ΔDOR points over the data arc, which is also believed to have contributed to the relatively poor performance that was obtained. This aspect of Earth-1 navigation operations is described in greater detail by Gray [6]. As is evident in both Figs. 4 and 5, the ranging data, when utilized at a 1-km weight, have little effect on the Doppler-plus-range solutions versus the Doppler-only solutions, indicating that at 1 kilometer most of the information content of the ranging data is being effectively discarded, except for the direct measurement of the Earth-spacecraft distance, r .

In the solutions obtained with a 2-mm/sec Doppler weight (Fig. 5), the Doppler-only and Doppler/1-km range solutions improved noticeably in terms of the error relative to the post-flyby reconstruction. Additionally, the dispersion ellipses for these two cases were more commensurate with the actual orbit determination errors than in the 1-mm/sec Doppler cases. This indicates that with a weight of 1 mm/sec, the Doppler data were being affected by some modeling error that was not adequately accounted for by either the assumed level of random measurement noise or with the estimated and consider parameter set described in Table 1. It is believed that the principal error source causing this behavior was the solar plasma effect (larger than expected ionospheric calibration errors may have also been a contributing factor), which was not explicitly accounted for by the Doppler measurement error model used in the

orbit determination filter, but was accounted for in the precision range error model by the stochastic bias parameters (hence, the good agreement between the aiming plane error for the 10-m range-only solution and the dispersion ellipse associated with this solution). As in the solution set with 1-mm/sec Doppler, the accuracy of the Doppler/1-km range/ ΔDOR solution with 2-mm/sec Doppler was found to be poor ($0.43 \Delta\text{rad}$) relative to the theoretical angle-finding capability of the ΔDOR data ($0.04\text{--}0.08 \mu\text{rad}$). It should be noted here that theoretical studies have indicated that navigation accuracies of $0.08\text{--}0.10 \mu\text{rad}$ may be achieved with two-way X-band (8.4-GHz) ranging and Doppler data, provided that accurate (1 m or better) station delay and transmission media calibrations are available [1].

VI. Conclusions

The results of the analysis indicate that for Galileo Earth-1 approach navigation, precision ranging data yielded orbit solutions which, although not in theory as accurate as those obtainable with ΔDOR data, were in fact somewhat better in this particular case. The relatively poor navigation performance of ΔDOR is primarily attributed to the sparsity of the ΔDOR data set and the irregular distribution of these measurements within the data arc. In addition, it was found that the orbit determination errors obtained in the Doppler and precision ranging solutions were consistent with simple theoretical predictions of the angle-finding capability of S-band ranging and Doppler data. The relative ease of ranging data scheduling and processing procedures makes ranging an attractive alternative to ΔDOR , when the nominally higher performance of ΔDOR is not required. Further improvements (factors of 3 to 5) in the navigation accuracy obtainable with precision ranging may be achieved through the use of X-band (8.4-GHz) frequencies, as opposed to S-band (2.3-GHz) frequencies, and through improved station delay and transmission media calibration accuracies.

Acknowledgments

The authors gratefully acknowledge the many valuable comments and suggestions received from H. W. Baugh, J. K. Campbell, G. S. Johnson, J. R. Smith, and A. H. Taylor. Special thanks are due to R. A. Jacobson, who helped develop some of the concepts that were ultimately used in this analysis, and to S. P. Synnott, for his careful review of the manuscript for this article.

References

- [1] S. W. Thurman, T. P. McElrath, and V. M. Pollmeier, "Short-Arc Orbit Determination Using Coherent X-Band Ranging Data," paper AAS-92-109, AAS/AIAA Spaceflight Mechanics Meeting, Colorado Springs, Colorado, February 24-26, 1992.
- [2] V. M. Pollmeier and P. H. Kallemeyn, "Galileo Orbit Determination from Launch Through the First Earth Flyby," *Proceedings of the 47th Annual Meeting of the Institute of Navigation*, Williamsburg, Virginia, pp. 9-16, June 10-12, 1991.
- [3] T. W. Hamilton and W. G. Melbourne, "Information Content of a Single Pass of Doppler Data from a Distant Spacecraft," *JPL Space Programs Summary 37-39*, vol. 3, pp. 18-23, March-April 1966.
- [4] J. D. Anderson, "The Introduction of Range Data Into the Data Compression Scheme," *JPL Space Programs Summary 37-43*, vol. 3, pp. 18-24, November-December 1966.
- [5] A. H. Taylor, J. K. Campbell, R. A. Jacobsen, B. Moultrie, R. A. Nichols, Jr., and J. E. Riedel, "Performance of Differenced Range Data Types in Voyager Navigation," *Journal of Guidance, Control, and Dynamics*, vol. 7, no. 3, pp. 301-306, May-June 1984.
- [6] D. L. Gray, " Δ VLBI Data Performance in the Galileo Spacecraft Earth Flyby of December 1990," *The TDA Progress Report 42-106*, vol. April-June 1991, pp. 335-352, August 15, 1991.
- [7] L. Efron and R. J. Lisowski, "Charged Particle Effects to Radio Ranging and Doppler Tracking Signals in a Radially Outflowing Solar Wind," *JPL Space Programs Summary 37-56*, vol. 2, pp. 61-69, January-February 1969.

Table 1. Galileo orbit determination model assumptions.

Model parameters	A priori uncertainty, 1σ	Remarks
<u>Estimated</u>		
Spacecraft state vector		No information
Epoch position	10^8 km	
Velocity	10^8 km/sec	
Solar radiation pressure		
Radial	5 percent of nominal	
Transverse	1 percent of nominal	
Attitude update maneuvers	0.5 mm/sec	About 1 every 2 weeks
Propellant line flushings		
Magnitude	0.5 mm/sec	About 1 every 3 weeks
Direction	15 mrad	
Quasar location, for Δ DOR	100 nrad	Conservative
Range bias parameters (1 per station-pass)		
DSN 70 m	2.0 m	
DSN 34 m STD	10.0 m	Conservative
<u>Consider</u>		
DSN station locations (correlated covariance), m		
Spin radius	0.24	Relative uncertainty between stations is approximately 5 cm
Longitude	0.24	
z-height	0.30	
Troposphere zenith delay calibration error, cm		
Wet	4.0	
Dry	1.0	
Ionosphere zenith delay calibration error, cm		S-band values (conservative)
Daytime	75.0	
Nighttime	15.0	
Acceleration biases, km/sec ²		
Radial (spacecraft spin axis)	3×10^{-12}	
Transverse	1×10^{-12}	
Earth ephemeris (heliocentric), km		A priori covariance, JPL ephemeris DE 125
Radial	0.2	
Along track	30.0	
Out-of-plane	15.0	
Earth mass, GM	$0.15 \text{ km}^3/\text{sec}^2$	DE 125

Table 2. Comparison of orbit solutions and reconstructed trajectory.

Case	Aiming plane error, km	Equivalent angular error, μrad
1-mm/sec Doppler (no range) (no ΔDOR)	27.7	2.22
1-mm/sec Doppler, 1-km range (no ΔDOR)	27.8	2.22
1-mm/sec Doppler, 10-m range (no ΔDOR)	4.2 ^a	0.34 ^a
1-mm/sec Doppler, 1-km range, 50-cm ΔDOR	7.8	0.62
2-mm/sec Doppler (no range) (no ΔDOR)	15.1	1.21
2-mm/sec Doppler, 1-km range (no ΔDOR)	12.1	0.97
2-mm/sec Doppler, 10-m range (no ΔDOR)	6.5 ^a	0.52 ^a
2-mm/sec Doppler, 1-km range, 50-cm ΔDOR	5.4	0.43
10-m range (no Doppler) (no ΔDOR)	7.3 ^a	0.58 ^a

^aPrecision ranging solutions.

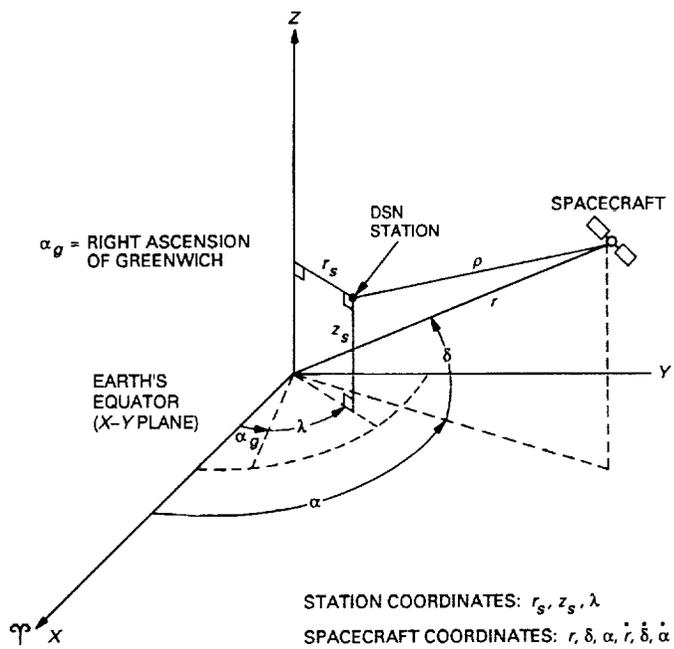


Fig. 1. Station-spacecraft tracking geometry.

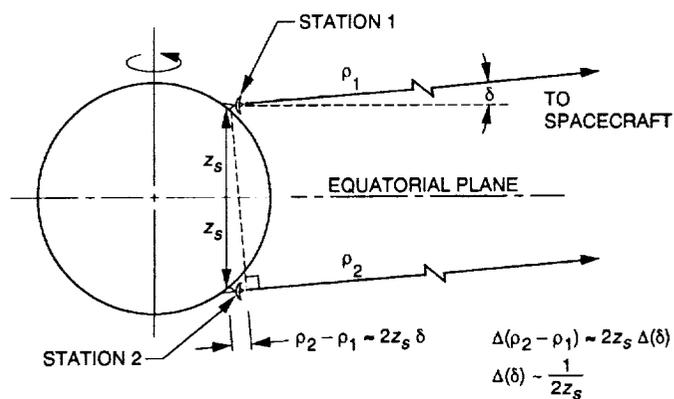


Fig. 2. S-band ranging and Doppler theoretical angular precision.

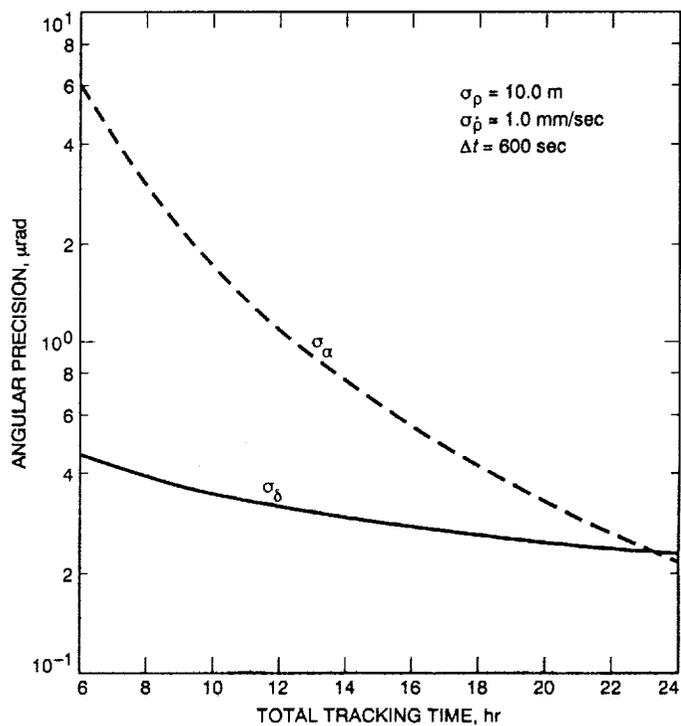


Fig. 3. Declination determination using range measurements from two widely separated stations.

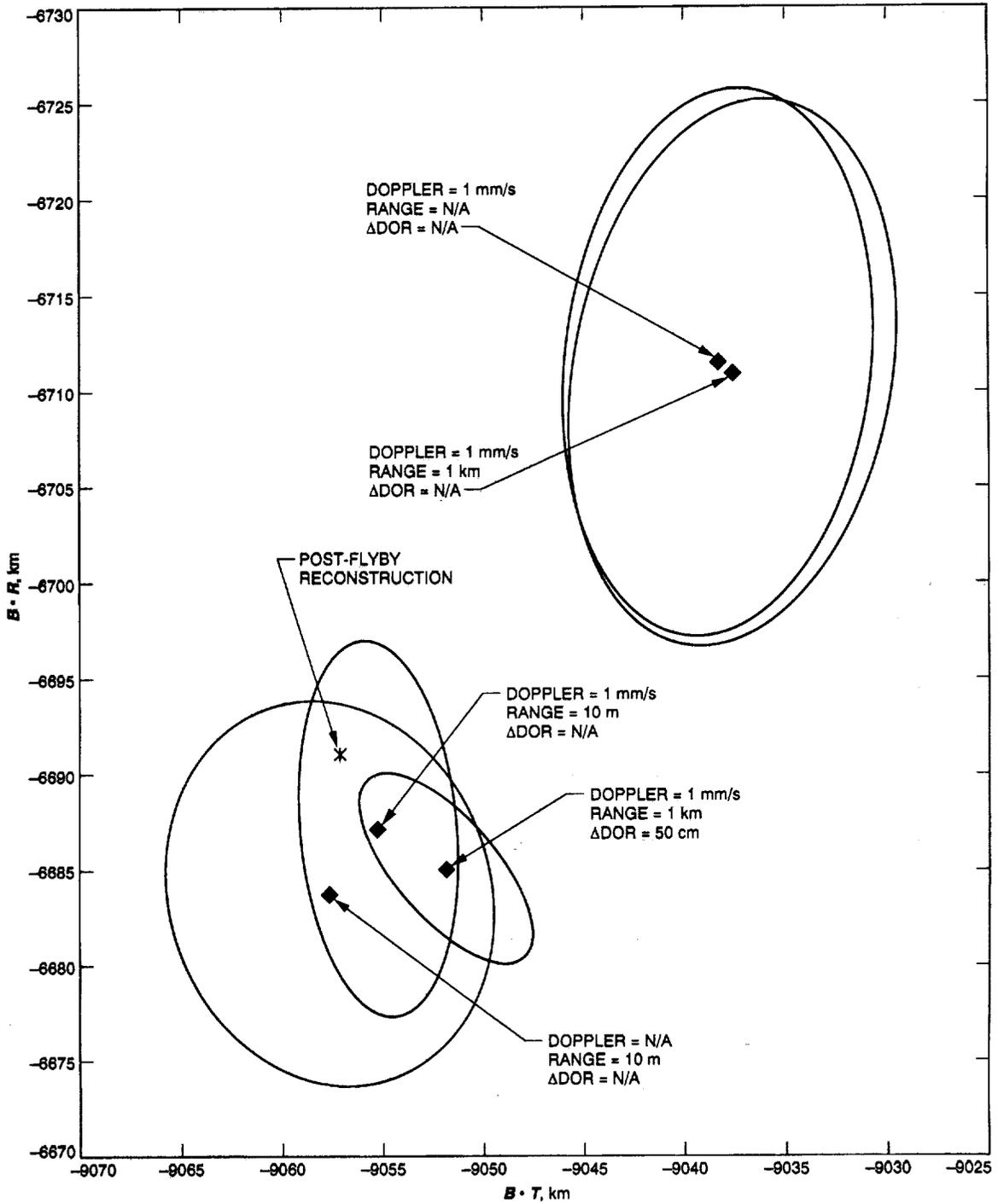


Fig. 4. Earth-1 aiming plane (solutions with 1-mm/sec Doppler weight).

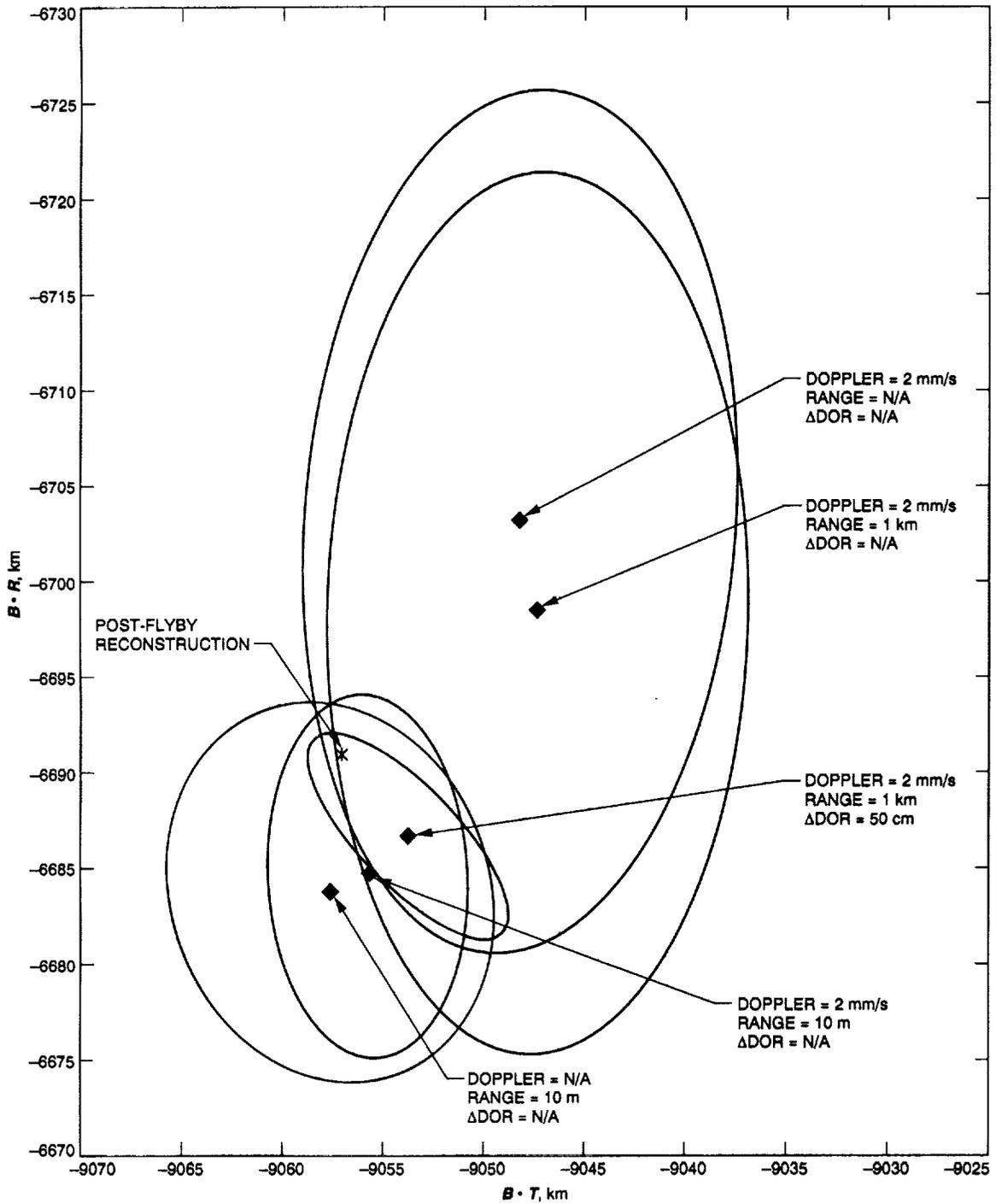


Fig. 5. Earth-1 aiming plane (solutions with 2-mm/sec Doppler weight).

N93-19416

53-46

128436

p-19

The Effect of Tropospheric Fluctuations on the Accuracy of Water Vapor Radiometry

J. Z. Wilcox

Tracking Systems and Applications Section

Line-of-sight path delay calibration accuracies of 1 mm are needed to improve both angular and Doppler tracking capabilities. Fluctuations in the refractivity of tropospheric water vapor limit the present accuracies to about 1 nrad for the angular position and to a delay rate of 3×10^{-13} sec/sec over a 100-sec time interval for Doppler tracking. This article describes progress in evaluating the limitations of the technique of water vapor radiometry at the 1-mm level. The two effects evaluated here are: (1) errors arising from tip-curve calibration of WVR's in the presence of tropospheric fluctuations and (2) errors due to the use of nonzero beamwidths for water vapor radiometer (WVR) horns. The error caused by tropospheric water vapor fluctuations during instrument calibration from a single tip curve is 0.26 percent in the estimated gain for a tip-curve duration of several minutes or less. This gain error causes a 3-mm bias and a 1-mm scale factor error in the estimated path delay at a 10-deg elevation per 1 g/cm^2 of zenith water vapor column density present in the troposphere during the astrometric observation. The error caused by WVR beam averaging of tropospheric fluctuations is 3 mm at a 10-deg elevation per 1 g/cm^2 of zenith water vapor (and is proportionally higher for higher water vapor content) for current WVR beamwidths (full width at half maximum of approximately 6 deg). This is a stochastic error (which cannot be calibrated) and which can be reduced to about half of its instantaneous value by time averaging the radio signal over several minutes. The results presented here suggest two improvements to WVR design: First, the gain of the instruments should be stabilized to 4 parts in 10^4 over a calibration period lasting 5 hours, and second, the WVR antenna beamwidth should be reduced to about 0.2 deg. This will reduce the error induced by water vapor fluctuations in the estimated path delays to less than 1 mm for the elevation range from zenith to 6 deg for most observation weather conditions.

I. Introduction

Future missions will benefit from spacecraft tracking with 100-prad accuracy in an angular position. The Cassini radio science team requires a delay-rate accuracy of 5×10^{-16} sec/sec over 1000 sec for gravitational wave searches using Doppler tracking. Fluctuations in wet tropospheric refractivity limit the present tracking capabilities to about 1 nrad for angular position and 5×10^{-14} sec/sec over 1000 sec for Doppler tracking, which corresponds to a 0.1-mm/sec uncertainty in spacecraft radial velocity.

Angular tracking is done with very long baseline interferometry (VLBI), a technique that measures the differential phase between two DSN antennas of an electromagnetic signal originating from a radio source. By relating this measured phase difference to geometrical path delays, astrometric parameters can be estimated. Inhomogeneities in tropospheric refractivity cause unmodeled errors in the path delay at about the 1-cm level for a path delay at zenith over a period of several hours. This error corrupts angular position estimates at about the 1-nrad level. To achieve 100-prad angular position accuracy, fluctuations in path delays must be calibrated at the 1-mm level [1,2].

Charged particles (both in the Earth's ionosphere and in solar wind) are the main source of error for spacecraft gravitational wave searches using Doppler tracking at S-band (2.3 GHz) (all missions prior to Galileo). Searches planned for Galileo at X-band (8.4 GHz) will be limited by fluctuations in the plasma and the troposphere at about the same level, a delay rate of approximately 5×10^{-15} sec/sec over 1000 sec [3]. In searches planned for Cassini at K-band (32 GHz), the increased observational frequency will reduce the plasma-induced error to approximately 5×10^{-16} sec/sec over 1000 sec. To take advantage of this increased sensitivity, tropospheric fluctuations must be calibrated at the submillimeter level.

Water vapor radiometers (WVR's) have been suggested for measuring line-of-sight path delays due to tropospheric water vapor. WVR's work on the principle that the radio power collected by a WVR antenna is proportional to the brightness temperature of the sky in the antenna pointing direction. Path-delay retrieval algorithms relate the brightness temperatures measured at two (or more) frequencies near the 22.6-GHz water vapor absorption line to a path delay associated with tropospheric water vapor along a line of sight in the same direction [4]. Brightness temperatures at two or more frequencies are required to

subtract the effect of liquid water absorption from the total absorption. Unlike water vapor, liquid water does not affect the refractivity at microwave frequencies.

The ability of WVR's to determine the absolute path delays, or to track path delay changes induced by tropospheric fluctuations, has been tested recently with mixed results.¹⁻³ The accuracy of WVR calibration using radiosondes is limited at a level of approximately 10 percent (which causes a 0.6-cm error for a 6-cm zenith path delay) due to uncertainties in both the radiosonde data and path delay retrieval algorithms [4].² Studies that calibrated VLBI time series with differenced WVR delays using co-pointed antennas reported reduced rms residual delays at high elevations,³ but the residual delays actually increased (up to 20 psec, corresponding to a 0.6-cm path delay error) at elevations below 50 deg. These results strongly suggest that to meet the future mission requirements, the character of various error sources that affect the accuracy of water vapor radiometry must be much better understood.

The accuracy with which WVR's estimate tropospheric path delays is determined by error sources that include (but are not restricted to) inaccuracies in calibration of the WVR gain, uncertainties in the path-delay retrieval algorithms and atmospheric absorption modeling, radiometer noise, the effect of the WVR location relative to the radio telescope, and the beam intensity distribution. The error sources will be discussed later. The errors can originate in the instrument or in the atmosphere and can cause bias, scale, or random errors in the retrieved path delays. The effect of a constant bias on path delays for angular tracking can be eliminated by differencing between observations. Biases have no effect on delay rates used for gravitational wave searches. Unmodeled variability in spatial distributions of atmospheric parameters together with the system noise have been recognized as the principal noise mechanisms that limit the ability of WVR's to monitor tropospheric fluctuations. Only recently has it been recognized that for any realistic WVR design, the tropospheric dynamics will also influence the WVR accuracy.

¹ T. J. Vesperini, "Stapleton WVR Experiment-Part I: Results," JPL Interoffice Memorandum (internal document), Jet Propulsion Laboratory, Pasadena, California, March 7, 1988.

² S. Keihm, "Water Vapor Radiometer Intercomparison Experiment: Platteville, Colorado, March 1-14, 1991," *Final report, JPL Task Plan 80-3289* (internal document), Jet Propulsion Laboratory, Pasadena, California, July 1991.

³ C. Edwards, "Water Vapor Radiometer Line-of-Sight Calibration Capabilities," JPL Interoffice Memorandum 335.1-90-015 (internal document), Jet Propulsion Laboratory, Pasadena, California, March 30, 1990.

The main goal of this article is to investigate how the tropospheric dynamics affect the WVR's ability to track tropospheric fluctuations. Because WVR's are imperfect instruments, one of the goals of this article is to quantify relations between WVR design parameters and path-delay retrieval accuracy. To understand how the errors affect path delay estimates, an analysis of the effect of WVR error sources on path delay retrieval is presented. The article then focuses on two specific errors that arise from tropospheric fluctuations: errors in WVR instrument gain calibration from tip curves and errors caused by WVR antenna beam averaging. Tropospheric fluctuations cause errors in the estimated gain because data at different tip directions are analyzed using mapping relationships valid for a temporally constant and spatially homogeneous troposphere. Tropospheric fluctuations cause unmodeled departures from this picture, which induces errors in the estimated gain and ultimately in the estimated path delays. WVR antenna beam averaging causes errors in the estimated brightness temperature in the direction of a radio telescope because data recorded by the WVR's are beam averaged around the WVR pointing direction (full beam widths of the present state-of-the-art WVR's are between 4 and 10 deg), whereas radio telescopes measure the tropospheric effects along the line of sight to a distant radio source. The averaging corrupts the accuracy of the estimated brightness temperature in the radio source direction for all realistic WVR designs.

The article is organized as follows: Section II discusses how various error sources affect path delay estimates. Some basic equations describing the conversion of the recorded data to the line-of-sight path delay are given. Section III calculates the error in the WVR measurement incurred by using tip curves in an inhomogeneous troposphere. This error will exist as long as the WVR's are calibrated by using tip curves. Section IV studies the effect of beam averaging on WVR measurements for collocated antennas. The effects of antenna copointing and beam intensity distribution are discussed. The aim here was not to derive exact numbers for any specific WVR design but rather to provide rough error estimates. Section V is a summary with recommendations for WVR gain stability requirements and antenna beam width design to comply with the 1-mm path-delay accuracy requirement, as well as plans for further studies.

II. WVR Error Sources and Their Effect on Path Delay Retrieval

To put the above-mentioned effects in perspective, this section gives an overview of how various error sources affect the estimated path delay. The conclusions of this

overview are summarized in Fig. 1. Error sources result from inaccuracies in WVR measurements or data interpretation and can cause bias, scale, or random effects in the retrieved path delay. The effect of a constant bias or linear trend on the tracking of changes in the path delay can be eliminated with astrometric parameter estimation. A scale error is an error that is directly proportional to the path delay. It is unclear at the present time to what extent radiometric data reduction can filter out the effect of a systematic scale error on astrometric parameter estimates. Stochastic errors cannot be reduced by parameter estimation. The effect of error sources on the estimated path delays can be analyzed by using a path delay retrieval algorithm that relates the line-of-sight path delay (L_v) due to water vapor in the troposphere to the brightness temperature (T_B) in the same direction. For example, Eq. (20) of [4],

$$L_v = a_0 + a_1 T_{B,1} + a_2 T_{B,2} \quad (1)$$

expresses L_v as a linear combination of the brightness temperatures $T_{B,j}$ ($j = 1,2$) at frequencies $f_1 = 20.6$ GHz and $f_2 = 31.4$ GHz. In Eq. (1), a_j 's are coefficients that have been determined by a regression analysis of the WVR data [4]. The value of a_0 quantifies the effect of the dry component of the atmospheric absorption; a_1 quantifies the effect of water vapor absorption; and the term $a_2 T_{B,2}$ ($a_2 \simeq -a_1 (f_1/f_2)^2$) subtracts the contribution due to tropospheric liquid water from the total absorption (liquid water causes a negligible path delay). For L_v in centimeters and $T_{B,j}$ in kelvins, typical regression coefficients are $a_0 \simeq -1.6$, $a_1 \simeq 0.66$, and $a_2 \simeq -0.3$. Typical $T_{B,j}$'s are between 10 and 100 K. The value of L_v scales as $L_v(\text{cm}) \simeq 6 N_v / \sin E$, where N_v is the water vapor column density (in g/cm^2) at zenith, and E is the observed elevation. For typical N_v between 1 g/cm^2 and 4 g/cm^2 , the zenith path delay ($L_{v,z}$) is between 6 and 24 cm.

Taking the differential of Eq. (1) determines the error in the retrieved path delay, δL_v , in terms of errors in the regression coefficients, δa_j 's, and the errors in $T_{B,j}$'s as

$$\delta L_v = \delta a_0 + \sum_{j=1}^2 (\delta a_j T_{B,j} + a_j \delta T_{B,j}) \quad (2)$$

where the subindex $j = 1,2$ refers to the WVR frequency channels. This article focuses primarily on errors in $T_{B,j}$ induced by tropospheric fluctuations. However, for the completeness of the discussion and because they are so large, sources of errors in a_j 's are briefly summarized first.

The biggest errors in a_j 's come from two sources: (1) inaccurate modeling of the absorptivity of water vapor, α_v [i.e., from errors of the dependence of $\alpha_v = \alpha_v(T, p, \rho, f)$ on the ambient temperature T , pressure p , water vapor density ρ , and frequency f], and (2) uncertainties in the spatial and temporal distributions of atmospheric pressure, temperature, and water vapor and liquid present in the troposphere during the observation along the line of sight. Errors in the a_j 's cause scale errors in L_v , of the type $\delta L_v \simeq L_v \delta a_1 / a_1$. The error caused by inaccurate modeling of $\alpha_v(T, p, \rho, f)$ is systematic, about 10 percent for the current absorptivity models. Atmospheric profile uncertainties depend on season, site, weather, time of day, and line of sight. The uncertainties are difficult to model and are the reason why there is no one-to-one correspondence between the brightness temperature and the path delay. The error caused by uncertainties in atmospheric profiles is between 2 and 4 percent, depending upon the specific retrieval algorithm used. Several possibilities pertaining to the feasibility of reducing the δa_j 's will be noted in Section V.

The error in the brightness temperature, $\delta T_{B,j}$, comes from two sources: (1) the error $\delta T_{A,j}$ in the measured WVR antenna temperature $T_{A,j}$, and (2) an error in the interpretation of $T_{A,j}$. The latter contribution to $\delta T_{B,j}$ originates in the fact that WVR's do not measure $T_{B,j}$ directly, but rather they record a signal, $V_{A,j}$ (from which $T_{A,j}$ is extracted), and the obtained $T_{A,j}$ is used to estimate the brightness temperature $T_{B,j}$. The error $\delta T_{A,j}$ in the WVR-measured $T_{A,j}$ can be expressed in terms of the WVR parameters, as follows: All WVR's use an internal reference, such as blackbody radiation or a noise diode, to enable a subtraction of the contribution of the system temperature from the recorded data. The recorded $V_{A,j}$ is proportional to the difference between $T_{A,j}$ and the reference load temperature, T_{ref} ,

$$V_{A,j} = g (T_{A,j} - T_{ref}) \quad (3)$$

where g is the WVR gain. Inverting Eq. (3) to obtain $T_{A,j}$ and taking a differential of the resulting equation yields the error $\delta T_{A,j}$ in the measured $T_{A,j}$

$$\begin{aligned} \delta T_{A,j} &= \frac{\delta V_{A,j}}{g} + \delta T_{ref} - \frac{\delta g}{g^2} V_{A,j} \\ &= \frac{\delta V_{A,j}}{g} + \delta T_{ref} + \frac{\delta g}{g} (T_{ref} - T_{A,j}) \quad (4) \end{aligned}$$

Thus, $\delta T_{A,j}$ comes from three WVR parameters, the system noise (modeled as the uncertainty in the equivalent temperature, $\delta V_{A,j}/g$), and the uncertainties in the refer-

ence load temperature (δT_{ref}) and the WVR gain (δg). Figure 1 shows that for the path delay error to be less than 1 mm, the noise (or more generally, any unmonitored drifts) in all temperature-like quantities must be less than 0.2 K. In practice, the effect of the system noise on path delay estimates can be reduced (at the expense of time resolution) by increasing the signal integration time. As long as T_{ref} remains constant, the path delay error caused by using an incorrect value of T_{ref} is a constant bias (whose effect on astrometric estimates can be eliminated by differencing between the observations).

Because the antenna voltage $V_{A,j} \propto T_{A,j} - T_{ref}$, the error caused by δg [the last term on the right-hand side of Eq. (4)] consists of two terms. The first is $\propto T_{ref}$, the second is $\propto T_{A,j}$. Note that since the typical $T_{ref} \simeq 300$ K is bigger than $T_{A,j}$ ($T_{A,j} \simeq T_{B,j} \simeq 10$ to 100 K), the presence of T_{ref} in $V_{A,j}$ enhances the effect of the gain error (whatever its origin may be) on δL_v . By using Eq. (4) in Eq. (2), the first term (i.e., the term $\propto T_{ref}$) leads to $\delta L_v \propto a_i T_{ref} \delta g/g$. For stable gain (i.e., for a constant difference between the estimated and the WVR true gains), this is a constant bias. For unmonitored gain fluctuations, this is a random error, which, in order to satisfy the 1-mm path delay requirement (see Fig. 1), must be $\delta g/g < 0.08$ percent. The second term (i.e., the term $\propto T_{A,j}$) is a scale error, $\delta L_v \simeq L_v \delta g/g$, which is systematic for a stable gain and time varying for an unstable gain. This error depends on the tropospheric humidity and the observed elevation. For the error to be less than 1 mm, the gain error must be $\delta g/g < 0.1 \sin E / L_{v,z}$ (where $L_{v,z}$ is the zenith path delay in centimeters and $L_{v,z} = 6$ cm for a troposphere with 1 g/cm² of water vapor column density at zenith). Figure 1 shows that for the scale error to be less than 1 mm at a 10-deg elevation when $L_{v,z} = 6$ cm, $\delta g/g$ must be less than 0.26 percent. For the error to be less than 1 mm at $E = 6$ deg when $L_{v,z} = 24$ cm, $\delta g/g$ must be less than 0.04 percent.

The gain error can originate from several sources. The most troublesome of these are unmonitored instrumental drifts on a time scale of 100 to 1000 sec^{4,5} (stochastic fluctuations on time scales much shorter than the radio observations can be incorporated into the system noise). However, since they originate in the instrument, the drifts should be controllable by improved instrument stabiliza-

⁴ G. M. Resch (Tracking Systems and Applications Section) and S. Keihm (Microwave Observational Systems Section), personal communication, Jet Propulsion Laboratory, Pasadena, California, 1981.

⁵ G. Parks, C. Ruf, and S. Keihm, "Advanced Water Vapor Radiometer: Definition Phase Study" (internal document), Jet Propulsion Laboratory, Pasadena, California, November 29, 1990.

tion in advanced WVR designs. Tropospheric fluctuations are another source of δg for WVR's calibrated using the tip curves. Even though this δg is stochastic in origin, the ensuing difference between the estimated and the actual WVR gains causes a systematic (bias and scale) error in L_v . The error depends on tropospheric humidity and calibration strategy and will be calculated in Section III of this article.

The other contribution to $\delta T_{B,j}$ is the error made in inferring $T_{B,j}$ from $T_{A,j}$. The simplest and most often used relationship between $T_{B,j}$ and $T_{A,j}$ is that $T_{B,j} = T_{A,j}$. This relationship can be in error because of inaccurate antenna pointing and spatial separation between the WVR antenna and the radio telescope. Another source of uncertainty is the effect of averaging tropospheric fluctuations over WVR beam intensity distribution. Section IV derives an expression for the error in the estimated $T_{B,j}$ that was caused by the beam averaging of tropospheric fluctuations for collocated and copointed antennas. The error will increase with decreasing elevation more rapidly than L_v , i.e., faster than a simple scale error. Since, as a result of this increase, low-elevation data will be weighted more heavily than they should be during VLBI data reduction, astrometric parameter estimates will be impacted.

III. Error in the WVR Gain Estimated From Tip Curves Due to Tropospheric Fluctuations

This section presents a calculation of the error in the estimated gain (\hat{g}) of WVR's induced by tropospheric fluctuations during the WVR calibration using the tip curves.^{6,7} Tip curves use the elevation dependence of the sky brightness temperature [$T_{B,i} \equiv T_B(E_i)$, where E_i is the tip elevation] to calibrate the WVR gain. The WVR gain is determined by fitting the WVR recorded signal $V_{A,i}$ ($V_{A,i} \equiv V_A(E_i)$). If the elevation dependence of $T_{B,i}$, and therefore of $V_{A,i}$, were known, this calibration procedure would be limited only by thermal measurement noise. Spatial and temporal fluctuations of atmospheric water vapor cause the actual distribution of water vapor to depart from the static distribution assumed in fitting the $V_{A,i}$'s.

⁶ J. Z. Wilcox, "The Standard Deviation of WVR Gain Estimated from Tip Curves due to Wet Troposphere Fluctuations," JPL Interoffice Memorandum 335.6-91-032 (internal document), Jet Propulsion Laboratory, Pasadena, California, December 19, 1991.

⁷ J. Z. Wilcox, "The Difference Between Two Successive WVR Gain Estimates From Tip Curves due to Wet Troposphere Fluctuations," JPL Interoffice Memorandum 335.6-91-033 (internal document), Jet Propulsion Laboratory, Pasadena, California, December 20, 1991.

Therefore, an error due to tropospheric fluctuations is introduced into the gains estimated from the tip curve data. All delays calculated with the derived gain will, therefore, also be in error. In this section, the covariance of the gain estimates is determined in terms of the covariance of the tropospheric opacity. This is done by performing a tip-curve analysis of modeled data and evaluating the opacity covariance by using the Kolmogorov turbulence model [1] for the wet troposphere.

A. Tip-Curve Analysis

The tip-curve data were modeled by using the optically thin tropospheric approximation for the standard radiation transport equation, neglecting the effect of the Earth's curvature, ray bending, and nonzero WVR beamwidth, assuming negligible time elapsed during each tip curve sequence, and neglecting the time variation of all other model parameters except the tropospheric opacity. The recorded signal $V_{A,i}(t)$ for tip curve epoch t at elevation E_i is then expressed as [5]

$$\begin{aligned} V_{A,i}(t) &= g (T_{B,i}(t) - T_{ref}) \\ &= g (T_C e^{-\tau_i(t)} + T_M(1 - e^{-\tau_i(t)}) - T_{ref}) \\ &\simeq g (T_C + T_{MC} \tau_i(t) - T_{ref}) \end{aligned} \quad (5)$$

where the subscript i refers to i th elevation, g is the WVR gain, T_{ref} ($T_{ref} \simeq 300$ K) is the reference temperature discussed in Section II, $T_C \simeq 2.8$ K is the cosmic background temperature, $T_{MC} = T_M - T_C$ where $T_M \simeq 270$ to 280 K is the average atmospheric temperature [5], and $\tau_i(t)$ ($\tau_i(t) = \tau(E_i, t)$) is a time-varying line of sight opacity at elevation E_i . The linear approximation [the second expression on the right-hand side of Eq. (5)] is a good approximation for most optically thin ($\tau_i < 0.5$) tropospheres of interest, with the added benefit of mathematical simplicity.

In the standard tip curve analysis, the method of least squares [6] is used to obtain the gain estimate (\hat{g}) in terms of the tip data. Appendix A discusses the tip curve fitting in detail. The data are fit to a static (temporally averaged) version of Eq. (5) by using the mapping function $\langle \tau_i \rangle = \tau_z A_i$, where $\langle \dots \rangle$ designates the statistical average, τ_z is the averaged opacity mapped to zenith, and the air mass $A_i = 1/\sin E_i$. For N elevations, the fitting leads to an equation of the following type:

$$\hat{g}(t) = \sum_{i=1}^N c_i V_{A,i}(t) \quad (6)$$

where c_i 's are coefficients that depend on the assumed T_C , T_{MC} , and T_{ref} , and on the so-called variance-covariance matrix $W^{-1} \equiv \text{cov}(V_A(t), V_A(t'))$ (See [6] and Appendix A). Taking Eq. (6) at t and t' and substituting Eq. (5), the covariance of the gain estimates separated by the time interval $T = t - t'$ is obtained in terms of opacity correlations as

$$\begin{aligned} \text{cov}(\hat{g}(t), \hat{g}(t')) &= \sum_{i,j=1}^N c_i c_j \text{cov}(V_{A,i}(t) V_{A,i}(t')) \\ &\simeq g^2 T_{MC}^2 \sum_{i,j=1}^N c_i c_j \text{cov}(\tau_i(t) \tau_j(t')) \end{aligned} \quad (7)$$

The opacity covariance $\text{cov}(\tau_i(t) \tau_j(t'))$ was evaluated by neglecting fluctuations in the dry component of $\tau_i(t)$ (dry troposphere contributes less than about 30 percent of the total opacity fluctuations)⁸ and by describing the fluctuations in the wet component by Kolmogorov turbulence [1]. Specifically, the wet contribution to τ_i was expressed as the line-of-sight integral $\tau_{v,i} = \int_0^{A_i h_v} \alpha_v(\vec{r}_i, t) dr_i$, where $\alpha_v(\vec{r}, t)$ is the tropospheric absorptivity due to water vapor per unit length at \vec{r} , $dr_i = A_i dz$ is the path increment along the line of sight at elevation E_i , and h_v is the wet troposphere height. Using $\alpha_v(\vec{r}, t) \simeq \chi(\vec{r}, t) \tau_{v,z} / L_{v,z}$, where $\chi = \text{index of refraction} - 1$, and $\tau_{v,z}$ and $L_{v,z}$ are wet opacity and path delay at zenith, respectively, Appendix B shows explicitly the integral expressions that relate $\text{cov}(\tau_i(t) \tau_j(t'))$ to the structure functions for $\chi(\vec{r}, t)$. (Note that $L_{v,z} \simeq 6$ cm, and $\tau_{v,z} \simeq 0.04$ and 0.02 per 1 g/cm^2 of zenith water vapor column density, and the dry opacity $\tau_{d,z} \simeq 0.017$ and 0.04 , at 20.6 GHz and 31.4 GHz , respectively.) The refractivity structure functions were evaluated by generalizing the Kolmogorov turbulence expression [1] for the structure function $\langle (\chi(\vec{r}) - \chi(\vec{r} + \vec{R}))^2 \rangle$ to inhomogeneities correlated both spatially and temporally [1]:⁹

$$\begin{aligned} D_\chi(\vec{R}, T) &\equiv \left\langle (\chi(\vec{r}, t) - \chi(\vec{r} + \vec{R}, t + T))^2 \right\rangle \\ &= \frac{N_v^2 C^2 |\vec{R} + \vec{v} T|^{2/3}}{1 + (|\vec{R} + \vec{v} T| / L_s)^{2/3}} \end{aligned} \quad (8)$$

where \vec{R} and T are the spatial and temporal intervals over which the structure function is evaluated and \vec{v} is the wind velocity. The role of Eq. (8) in VLBI data reduction was

⁸ G. E. Lanyi, personal communication, Tracking Systems and Applications Section, Jet Propulsion Laboratory, Pasadena, California, 1991.

⁹ See Footnote 7.

discussed in [1]. By using the standard deviation for retrieved path delays at average DSN conditions, the turbulence strength was shown to be $C = 2.4 \times 10^{-7} \text{ m}^{-1/3}$ for a tropospheric slab with a water vapor column density of $N_v \simeq 1 \text{ g/cm}^2$ at zenith (corresponding to approximately 6 cm of wet path delay) and height $h_v = 1 \text{ km}$ [1]. For $h_v = 2 \text{ km}$, the recalculated $C = 1.1 \times 10^{-7} \text{ m}^{-1/3}$, which is the value used in this article.¹⁰ The turbulence saturation scale length was taken to be $L_s = 3,000 \text{ km}$ [1]. The temporal correlation depends on the wind velocity v (Reference [1] shows that if one identifies T with $|\vec{R}|/v$, the spatial correlation between two parallel lines of sight separated by the distance $|\vec{R}|$ projected on the Earth's surface is equal to the temporal correlation of a single line of sight at time t and later $t+T$.) This article used $v = 10 \text{ m/sec}$, which is a typical wind speed at the Goldstone DSN antenna site.

B. The Estimated Gain Error

The standard deviation of the estimated gain, $\sigma_{\hat{g}}$, is equal to the square root of the covariance given by Eq. (7) for $t = t'$. For $t \neq t'$, Eq. (7) gives the gain covariance as a function of the time interval $T = |t - t'|$ between two successive gain estimates, $\sigma_{\hat{g}}(T) \equiv (\text{cov}(\hat{g}(t), \hat{g}(t+T)))^{1/2}$. Before discussing the numerical results, note that Eq. (7) depends on the number (N) of tip elevations. Evaluating Eq. (7) for $N = 2, 3$, and 4 , the covariance was found to depend on the tip range and be nearly independent of the elevation distribution within the range.¹¹ That is, the covariance is determined by the least-correlated tropospheric inhomogeneities (i.e., by correlations between the lines of sight associated with the minimum, E_{min} , and maximum, E_{max} , tip elevations). In what follows, the errors will be shown versus E_{min} (with E_{max} at zenith) per 1 g/cm^2 of zenith column density of water vapor present in the troposphere during WVR calibration. Note that because they are proportional to N_v , the errors are minimized by calibration in dry (and stable) weather.

The value of $\sigma_{\hat{g}}$ is shown in Fig. 2 as a function of E_{min} . The error has a flat minimum of approximately 0.26 percent between approximately $E_{min} = 10 \text{ deg}$ and 30 deg . Below a 10-deg elevation, the error increases with decreasing E_{min} because the decorrelation between the tropospheric fluctuations associated with E_{min} and 90-deg lines of sight increases more rapidly than does the air mass difference. Above a 30-deg elevation, the opposite is true. It can be shown that as E_{min} approaches 90 deg, the error increases as $(90 - E_{min})^{-3/2}$ (which is slower than it

¹⁰ See Footnote 6.

¹¹ See Footnote 6.

would be if the fluctuations were completely random, in which case the error would increase as $(90 - E_{min})^{-2}$. Thus, to minimize the gain error, E_{min} should be between 10 and 30 deg.

The least-squares fit coefficients c_i 's depend on the variance-covariance matrix W^{-1} . In the so-called consider analysis, the observable errors are assumed to be uncorrelated, and W^{-1} is approximated by a unit diagonal [1]. Using the observable variance-covariance matrix W^{-1} minimizes the variance of the estimated gain, Eq. (7) (i.e., it minimizes $\sigma_{\hat{g}}$). The errors calculated by using the unit and observable variance-covariance matrix W^{-1} are shown as solid and broken line curves, respectively, in Fig. 2. Note that the two curves are practically the same in the minimum region and differ by 8 percent at most in the wings. This indicates that using the full covariance-variance matrix W^{-1} does not significantly improve the estimated gain accuracy.

Figure 3 shows the temporal development of $\sigma_{\hat{g}}(T)$. The value of $\sigma_{\hat{g}}(T)$ decreases at $T \ll T_{corr}$ approximately as $\sigma_{\hat{g}}(T) \simeq \sigma_{\hat{g}} \sqrt{1 - (2T/T_{corr})}$, where T_{corr} is the decorrelation time. (Note that in Fig. 3, $\sigma_{\hat{g}}(T_{corr}/2) \simeq \sigma_{\hat{g}}/2$.) The value of T_{corr} depends on v and E_{min} as $T_{corr} \simeq h_v / (v \tan E_{min})$. That is, T_{corr} is the time it takes for a "frozen" troposphere [1] to pass through the tip range between E_{min} and zenith. For $v = 10$ m/sec, and $E_{min} \simeq 30$ deg, T_{corr} is about 7 min. The single-gain estimate errors ($\sigma_{\hat{g}}$'s) become independent from each other when the time T between subsequent tip curves exceeds T_{corr} . Therefore, if one wishes to minimize the estimated gain error by tip curve repetition, the tip curves should be separated by a time interval greater than 7 min.

C. The Estimated Path Delay Error

It has already been discussed in Section II that a WVR gain error induces two types of errors in the estimated L_v : a bias and a scale error. Using Eq. (4) in Eq. (2) and designating the bias and scale errors as ΔL_v and δL_v , respectively, the two errors are

$$\begin{aligned} \Delta L_v(\text{cm}) &\simeq (a_1 + a_2) T_{ref} \frac{\sigma_{\hat{g}}}{g} \\ &\simeq 120 \frac{\sigma_{\hat{g}}}{g} \simeq 0.3 N_{v,cal} \end{aligned} \quad (9a)$$

$$\begin{aligned} \delta L_v(\text{cm}) &\simeq (a_1 T_{A,1} + a_2 T_{A,2}) \frac{\sigma_{\hat{g}}}{g} \\ &\simeq L_v \frac{\sigma_{\hat{g}}}{g} \simeq 0.016 \frac{N_{v,cal} N_{v,obs}}{\sin E} \end{aligned} \quad (9b)$$

where the gain estimates at 20.6 GHz and 31.4 GHz were assumed to be correlated (i.e., the gains in the two WVR channels were determined during the same tip sequence) and their magnitudes the same. $T_{ref} = 300$ K, $N_{v,cal}$ and $N_{v,obs}$ are the number of grams per cm^2 of the column density of water vapor at zenith during the WVR calibration and radio observation, respectively, E is the elevation of the observation, and the last expressions on the right-hand sides correspond to the minimum $\sigma_{\hat{g}}/g \simeq 0.26$ percent. The ΔL_v corresponding to $N_{v,cal} = 1$ g/cm^2 (i.e., to $L_{v,z} = 6$ cm) is shown in Fig. 2. The minimum ΔL_v is approximately 3 mm. Note that as long as the WVR gain remains constant, ΔL_v is also constant, which makes it possible to remove its effect on astrometric estimates and delay rates by differencing between observations. For unstable gain, the gain changes must be monitored (to achieve a 1-mm path delay accuracy) with 0.08 percent accuracy. The dependence of the scale error (strictly speaking, δL_v will be a pure scale error only in optically thin tropospheres) on elevation was shown in Fig. 1. The value of δL_v caused by a 0.26-percent gain error at a 10-deg elevation when $N_{v,obs} = 1$ g/cm^2 is approximately 1 mm; when $N_{v,obs} = 4$ g/cm^2 , the error is 4 mm. To reduce δL_v to 1 mm, the gain error will have to be reduced by using a different calibration technique, tip curve repetition, or a parameter estimation during the data analysis. How this can be accomplished for a stable WVR is discussed in Section V of this article.

A note should be made here on the dependence of ΔL_v and δL_v on T_{ref} . From Eq. (9a), it would seem that $\Delta L_v \propto T_{ref}$. However, by performing the least-squares analysis, one finds that the tropospheric fluctuation-induced $\sigma_{\hat{g}}$ is $\propto 1/T_{ref}$. This cancels the dependence of ΔL_v on T_{ref} and makes $\delta L_v \propto 1/T_{ref}$. [Note, however, that ΔL_v caused by an instrumental gain drift will be $\propto T_{ref}$, as given by the first expression on the right-hand side of Eq. (9a).]

IV. Error in the Estimated Brightness Due to the WVR Nonzero Beamwidth

Retrieval algorithms relate the line of sight L_v to the brightness temperature (T_B) in the same direction, whereas data recorded by the WVR's are beam averaged around the WVR pointing direction. For collocated and copointed WVR antennas and radio telescopes (the telescope points along the line of sight to the radio source), the copointing introduces two types of errors into the estimated T_B : a systematic error due to the nonlinear dependence of air mass on elevation and a random error due to WVR beam

averaging of tropospheric fluctuations.¹² The systematic error can be calculated for known beam intensity distributions (beam shapes) and water vapor content in the troposphere. Its effect on path delay estimates can be eliminated by pointing the WVR to a slightly higher elevation so that the radio source lies in the direction of the centroid of the distribution of WVR beam brightness. In this section, the direction of the brightness centroid is calculated by WVR beam averaging of the air mass [see Eq. (12) for the centroid definition], and the random error is determined by using the tropospheric opacity statistically as in Section III. It has been suggested that to simplify WVR beam steering for collocated antennas, radio telescope and WVR antennas should be copointed.¹³ After correcting the WVR data by subtracting from them the systematic error, the "corrected" data would be used to estimate path delays in the direction of the beam's geometrical center. Therefore, the random error has been evaluated also for this geometrical center pointing case (and found that it is bigger than the random error for the brightness centroid pointing, by an amount that depends on the WVR beamwidth and elevation).

Before presenting numerical results, the WVR beam intensity distribution must be specified. The systematic error for an assumed Gaussian beam with 7.5-deg full width at half maximum (FWHM) has been calculated previously.¹⁴ To simplify the computations, and since the aim of this article is to provide error estimates (rather than to tailor the errors to specific WVR beam designs), a beam is used whose cross section when viewed in the propagation direction is a square with sharp cutoffs for the beam intensity. Specifically, for a beam centered at elevation E_c and azimuth φ_c , the WVR antenna temperature, $T_A(E_c)$, is calculated by integrating the brightness temperature of the sky ($T_B(E)$) over the intensity distribution:

$$T_A(E_c) \equiv [T_B]_c \simeq \int \int B_c(E, \varphi) T_B(E) dE \cos E d\varphi \quad (10)$$

where $[..]_c$ signifies the WVR beam average around (E_c, φ_c) , E and φ are elevation and azimuth angles, respec-

¹² J. Z. Wilcox, "The Error in the Estimated Path Delay due to WVR Antenna Beam Width: Beam Averaged Air Mass and Wet Troposphere Fluctuations Effects," JPL Interoffice Memorandum 335.6-92-004 (internal document), Jet Propulsion Laboratory, Pasadena, California, January 31, 1992.

¹³ See Footnotes 4 and 5.

¹⁴ S. Robinson, "A Simple Analytic Correction for WVR Beam Width," JPL Interoffice Memorandum 335.4-530 (internal document), Jet Propulsion Laboratory, Pasadena, California, July 23, 1985.

tively, and $B_c(E, \varphi)$ is the WVR antenna beam radiation pattern

$$B_c(E, \varphi) \simeq \frac{1}{(2 \Delta_{1/2})^2} \dots$$

when $|E - E_c| \leq \Delta_{1/2}$ and $|\varphi - \varphi_c| \leq \frac{\Delta_{1/2}}{\cos E}$ (11)

and zero otherwise, and $\Delta_{1/2} = \text{FWHM}/2$ is the WVR beam half-width. Note that as long as $E_c > \Delta_{1/2}$, ground pickup is avoided for this WVR beam pattern. By comparing the numerical results, the systematic error calculated using beam intensity distribution with sharp cutoffs (for the same FWHM) is about 20–30 percent smaller than for the Gaussian beams. (It was also found that neglecting beam spreading in the azimuthal direction underestimates the random error by less than 10 percent.) The radio telescope beam was approximated by an infinitely narrow pencil beam. Since the errors are proportional to tropospheric water vapor, all shown errors are for 1 g/cm² of zenith water vapor column density.

A. Tropospheric Fluctuation-Induced Error for Brightness Centroid Pointing

The elevation E_b of the WVR beam brightness centroid is determined by requiring that the statistically averaged brightness temperature $\langle T_B(E_b) \rangle$ at the centroid elevation E_b be equal to the WVR antenna temperature $\langle T_A(E_c) \rangle$

$$\langle T_B(E_b) \rangle = \langle T_A(E_c) \rangle \quad (12)$$

where E_c is the elevation of the WVR beam geometrical center, and $\langle \dots \rangle$ designates the statistical average. Substituting $T_B(E)$ [Eq. (5)] into Eq. (10), neglecting the effect of ray bending and the Earth's curvature, and using the optically thin troposphere approximation, the statistically averaged antenna temperature is

$$\langle T_A(E_c) \rangle = T_C + T_{MC} \tau_z [A]_c \quad (13)$$

where the WVR beam averaged air mass $[A]_c$ is given by the same integral expression as Eq. (10) except that $T_B(E)$ is replaced by $A_E = 1/\sin E$. By also using Eq. (5) for $T_B(E_b)$ in Eq. (12), one obtains the result that in an optically thin troposphere, the brightness centroid coincides with the air mass centroid,

$$\frac{1}{\sin E_b} = \int \frac{B_c(E, \varphi)}{\sin E} dE \cos E d\varphi \quad (14)$$

The calculated difference (the offset) between E_c and E_b is shown in Fig. 4. The offset is always positive (i.e., the brightness centroid is tilted from the beam's geometrical center to a lower elevation), it increases with the beam width approximately as $\Delta_{1/2}^2$, and it has a very wide minimum in the elevation range around $E_c \simeq 52$ deg. The minimum occurs because the systematic error, and hence the offset, depends on a nonvanishing second derivative of A_E versus E . Specifically, for $E < 90 - \Delta_{1/2}$, the offset $E_c - E_b \simeq (\Delta_{1/2}^2/6)A_c''/A_c'$, where A_c' and A_c'' are the first and second derivatives of $A_E = 1/\sin E$ versus E at E_c . Since A_c''/A_c' has a local minimum at 52 deg, the offset has

also a local minimum at 52 deg. Note also that as E_c approaches 90 deg, the offset rapidly increases to $\Delta_{1/2}/\sqrt{3}$. For the current WVR $\Delta_{1/2} < 4$ deg, the offset is less than 0.2 deg in the elevation range between approximately 20 and 80 deg. Note that in an optically thick troposphere, the brightness centroid will differ somewhat from the air mass centroid. This is a consequence of the nonlinear dependence of T_B on the air mass, Eq. (5), in an optically thick atmosphere.

The error in the estimated $T_B(E_b)$ is the square root of the variance

$$\sigma_T^2(E_b) = \langle (T_A(E_c) - T_B(E_b))^2 \rangle = \left\langle \left(\int T_B(E) (B_c(E, \varphi) - \delta(E - E_b, \varphi - \varphi_b)) dE \cos E d\varphi \right)^2 \right\rangle \quad (15)$$

where $\delta(E - E_b, \varphi - \varphi_b)$ is the Dirac delta function centered at (E_b, φ_b) . Equation (15) was evaluated by expressing $T_B(E)$ using Eq. (5), and then evaluating the correlations between the tropospheric opacities using the Kolmogorov turbulence model, as described in the paragraph following Eq. (7). The corresponding path delay error ($\sigma_{L,v}(E_b)$) was obtained by substituting $\sigma_T(E_b)$ into Eq. (2), where $\sigma_T(E_b)$'s were identified with $\delta T_{B,j}$'s ($j = 1, 2$ designates 20.6- and 31.4-GHz frequency channels, respectively) for the two WVR frequencies. Figure 5(a) plots the path delay error versus E_c . Note that the error increases with decreasing elevations faster than L_v . The error is plotted versus $\Delta_{1/2}$ in Fig. 5(b). After a rapid rise near zero, the error increases sublinearly. At $\Delta_{1/2} \simeq 3$ deg and $E_c = 30$ deg, 20 deg, and 10 deg, the errors are 0.06 cm, 0.1 cm, and 0.3 cm, respectively. Advanced WVR's with narrow beamwidths were designed for observations at low elevations. For $\Delta_{1/2} = 1$ deg, the errors are approximately 0.23 cm and 0.5 cm at $E_c = 10$ deg and 6 deg, respectively. These results qualitatively agree with errors calculated for copointed beams.¹⁵

The error shown in Fig. 4 refers to instantaneous measurements. However, astrometric data are averaged over time intervals on the order of 1 to 2 minutes, which tends to average out the fluctuations. Assuming that the radio telescope and the WVR observe simultaneously and continuously during t_{int} , the time averaged $T_B(E)$ is

$$\bar{T}_B(E) = \frac{1}{t_{int}} \int_0^{t_{int}} dt T_B(E, t) \quad (16)$$

Using $\bar{T}_B(E)$ instead of $T_B(E)$ on the right-hand side of Eq. (15), the result is shown in Fig. 6. At $t_{int} < T_{1/2}$, the error decreases with a time constant $T_{1/2} \simeq 2 \Delta_{1/2} h_v / v \sin^2 E_c$, which is the time required for the moving troposphere to pass through the WVR beam cone (and thus erase the tropospheric differences between the beam averaged and line-of-sight opacities). The value of the time constant $T_{1/2}$ increases with decreasing elevation and increasing half-width. For $h_v = 2$ km, $v = 10$ m/sec, $\Delta_{1/2} = 3$ deg, and $E_c = 30$ deg and 10 deg, $T_{1/2}$ is approximately 1.5 min and 12 min, respectively. For $\Delta_{1/2} \simeq 0.1$ deg, $T_{1/2}$ is 3 sec and 25 sec, respectively. Obviously, so that the time resolution is not degraded, WVR integration should never be longer than radio telescope integration.

B. Tropospheric Fluctuation-Induced Error for Geometrical Center Pointing

The systematic and random errors for a copointed radio telescope and a WVR antenna are calculated next. The usual argument for why WVR data should be associated with the direction of the beam's the brightness centroid is that for a constant troposphere, $\langle T_A(E_c) \rangle = \langle T_B(E_b) \rangle$. If this were the only criterion, one could also correct the WVR data (i.e., $T_A(E_c)$) by subtracting from them the estimated value of the difference between $\langle T_A(E_c) \rangle$ and $\langle T_B(E_c) \rangle$ and identify this "corrected" data as the actual value of the brightness temperature $T_B(E_c)$ in the direc-

¹⁵ S. Keihm, "Finite Beam Effects on LOS Path Delay Decorrelation," (internal document), Jet Propulsion Laboratory, Pasadena, California, March 22, 1990.

tion of the beam's geometrical center. For a copointed radio telescope and a WVR antenna, the (systematic) difference $\Delta T_A(E_c)$, is determined as

$$\Delta T_A(E_c) \equiv \langle T_A(E_c) - T_B(E_c) \rangle = T_{MC} \tau_z \Delta A_c \quad (17)$$

where τ_z is the total (wet and dry) zenith opacity, and $\Delta A_c \equiv [A]_c - A(E_c)$ is the difference between the beam averaged and geometrical center air mass. Note that $\langle T_A(E_c) \rangle$ is bigger than $\langle T_B(E_c) \rangle$.

The value of $\Delta T_A(E_c)$ at 20.6 GHz (and the corresponding path delay error) is shown in Fig. 7. The error increases with beamwidth approximately as $\Delta_{1/2}^2$. Note that for the radiation pattern with a sharp cutoff for intensity distribution, this quadratic dependence on $\Delta_{1/2}$ can be derived analytically

$$\Delta A_c = \frac{1}{2 \Delta_{1/2}} \ln \frac{\tan(E_c + \Delta_{1/2})/2}{\tan(E_c - \Delta_{1/2})/2} - A_c$$

$$\begin{aligned} \sigma_T^2(E_c) &= \langle (T_A(E_c) - \Delta T_A(E_c) - T_B(E_c))^2 \rangle \\ &= \left\langle \left(\int (T_B(E) - \langle T_B(E) \rangle) (B_c(E, \varphi) - \delta(E - E_c, \varphi - \varphi_c)) dE \cos E d\varphi \right)^2 \right\rangle \end{aligned} \quad (19)$$

where $\delta(E - E_c, \varphi - \varphi_c)$ is the Dirac delta function centered at (E_c, φ_c) . Equation (19) was evaluated by using the same procedure as Eq. (15). The corresponding path delay error ($\sigma_{L,v}(E_c)$) is shown in Fig. 8(a) versus E_c and in Figure 8(b) versus $\Delta_{1/2}$. Similarly as for the systematic error [and for the stochastic error for the brightness centroid, $\sigma_{L,v}(E_b)$], $\sigma_{L,v}(E_c)$ increases with decreasing E_c more rapidly than L_v . Note that when the systematic error is smaller than $\sigma_{L,v}(E_b)$ (such as for $\Delta_{1/2} < 1$ deg at a 6-deg elevation, or for $\Delta_{1/2} < 1.5$ deg at a 10-deg elevation), $\sigma_{L,v}(E_c)$ is about the same as $\sigma_{L,v}(E_b)$; whereas when the systematic error is bigger than $\sigma_{L,v}(E_b)$, $\sigma_{L,v}(E_c)$ looks more like the systematic error (which increases $\propto \Delta_{1/2}^2$ and can become very large). That is, while the copointing will significantly increase the stochastic error for wide WVR beams, the increase will be small for narrow beams. It has also been found that when $\sigma_{L,v}(E_c) \simeq \sigma_{L,v}(E_b)$ (the narrow beam case), the $\sigma_{L,v}(E_c)$ for the integrated signal [as in Eq. (16)] decreases with t_{int} at about the same rate as does $\sigma_{L,v}(E_b)$; whereas when $\sigma_{L,v}(E_c) > \sigma_{L,v}(E_b)$

$$\simeq \frac{\Delta_{1/2}^2}{6 \sin E_c} \left(1 + \frac{2}{\tan^2 E_c} \right) \quad (18)$$

where the approximate equality on the right-hand side has been obtained for narrow line widths, $\Delta_{1/2} \ll E_c$. The more important feature to notice in Fig. 7 is that the error increases with decreasing E_c more rapidly than L_v , namely that $\Delta L_v \propto \Delta T_A(E_c) \propto \Delta A_c \propto (1 + 2/\tan^2 E_c)/\sin E_c$. This is the same type of increase as for the random error calculated in the preceding paragraphs, i.e., the errors caused by WVR beam averaging increase with decreasing elevation more rapidly than a simple scale error. At $\Delta_{1/2} \simeq 3$ deg and $E_c = 30$ deg, 20 deg, and 10 deg, the systematic path delay error is approximately 0.36 mm, 1.2 mm, and 10 mm, respectively. The error decreases with decreasing beam width. At $\Delta_{1/2} \simeq 1$ deg and $E_c = 10$ deg and 6 deg, the error is 2 mm and 6 mm, respectively. (At $\Delta_{1/2} \simeq 0.1$ deg, the error is 0.1 mm and 0.06 mm, respectively.)

Correcting the measured $T_A(E_c)$ by $\Delta T_A(E_c)$, the stochastic error in the inferred $T_B(E_c)$ is determined from

(the wide beam case), the decrease is significantly slower.¹⁶ Thus, when the systematic error is less than $\sigma_{L,v}(E_b)$ (the narrow beam case, $\Delta_{1/2} < 1$ deg for all $E > 6$ deg), the copointing will introduce a negligible error into the estimated path delays using WVR's. However, for beam sizes greater than 1 deg, the WVR's should be pointed at a slightly higher elevation than the radio telescope.

V. Discussion and Recommendations

The main goal of this article is to investigate how tropospheric dynamics affect the ability of realistic WVR's to track tropospheric fluctuations. Two effects were studied in detail: errors in WVR instrument gain calibration from tip curves and errors in the estimated brightness temperature caused by WVR beamwidth averaging. The errors can be used to derive WVR gain stability requirements

¹⁶ See Footnote 12.

and WVR antenna beamwidth that would make it possible to reduce the path delay error to the 1-mm level in the elevation range from zenith to 6 deg.

The minimum error induced by tropospheric fluctuations in a single gain estimate (per 1 g/cm² of zenith water vapor column density during WVR calibration) is approximately 0.26 percent. That error causes two types of errors in the estimated path delay. The first error, approximately 3 mm, is independent of path delay. Provided that the WVR gain remains constant, this is a bias error that can be removed by differencing between VLBI observations (biases have no effect on delay rates used for gravitational wave searches). The second error, $\delta L_v \simeq L_v \sigma_{\dot{g}}/g$, is a scale error. Figure 9 shows the scale error for a single gain estimate δL_v (mm) $\simeq 0.16 N_{v,cal} N_{v,obs} / \sin E$ (where $N_{v,cal}$ and $N_{v,obs}$ come from water vapor content during the WVR calibration and radiometric observation) as a function of $N_{v,cal} N_{v,obs}$ and elevation (E) of the radiometric observation. For example, assuming that $N_{v,cal} = 1$ g/cm² (corresponding to $L_{v,z} = 6$ cm) and $N_{v,obs} = 2$ g/cm², δL_v exceeds 3 mm (i.e., it exceeds the L_v -independent error) when $E \leq 6$ deg. When $N_{v,cal} = N_{v,obs} = 2$ g/cm², δL_v exceeds 6.2 mm at $E = 6$ deg. To achieve the desired 1-mm path delay accuracy, either the effect of δL_v on astrometric estimates or δg itself must be reduced. The success of any approach to obtaining accurate astrometric estimates depends on the stability of the WVR gain.

For a stable WVR gain, the scale error is systematic. Preliminary results of attempts to reduce the effect of systematic scale errors on astrometric estimates with VLBI data analysis appear to be promising, although more work is needed to ascertain quantitative results.¹⁷ The gain error can be reduced by using an alternate gain calibration technique (such as two absolute reference load calibrations), or, assuming that the error induced by tropospheric fluctuations is the dominant error source, by tip curve repetition. Uncertainties in alternate calibration methods have so far prevented circumvention of tip curves. To reduce the error by tip curve repetition, the WVR gain must be sufficiently stable. For example, to achieve a 1-mm path delay accuracy at a 6-deg elevation when $N_{v,obs} = 2$ g/cm², the gain error must be less than 0.08 percent. To reduce the gain error to 0.08 percent when $N_{v,cal} = 2$ g/cm², the tip curve must be repeated at least $(0.52/0.08)^2 = 40$ times. Section III showed that successive gain estimates become decorrelated within a typ-

ical T_{corr} of approximately 7 min. Therefore, if one attempts to reduce the gain error by tip curve repetition, the tip curves should be separated by a time interval greater than 7 min, and the WVR gain should change by no more than 0.08 percent over at least $40 T_{corr} \simeq 5$ hr.¹⁸ (Because of the effect of other error sources and the possibility of humidity higher than 2 g/cm², the recommended WVR stability is 0.04 percent over the 5-hour period.) In the time interval between the WVR calibration and the radio metric observations (and during the observations) the gain will still have to be updated, e.g., by comparing the number of counts for the reference load.^{19,20}

The error in the estimated brightness temperature due to WVR beam averaging of tropospheric fluctuations was found to be smaller when the direction of the radio source coincided with the WVR beam brightness centroid than with the beam geometrical center. The error increases with tropospheric water vapor content and beamwidth. More important, however, is that the error increases with decreasing elevation faster than L_v . Since low-elevation data will be weighted more heavily than they should be, this will affect astrometric estimates. Advanced WVR designs have been suggested to reduce the error and avoid ground pickup by implementing narrow beams. Figure 10 shows the errors in the two-dimensional space of beam half-widths ($\Delta_{1/2}$) and elevations for 1 g/cm² of zenith water vapor column density. For each curve, the error is less than the cutoff error for all $\Delta_{1/2}$'s and E 's below and to the right of the curve. For example, for the error to be less than 1 mm at all $E > 10$ deg, $\Delta_{1/2}$ should be < 0.1 deg. For more humid weather, the errors will be higher (and the beamwidth requirement more stringent), proportional to zenith water vapor. Because of various approximations involved in deriving Fig. 10 (the WVR beam radiation pattern with sharp cutoffs for intensity distribution, an infinitely narrow pencil beam for the radio telescope, and an optically thin troposphere), the guidelines are approximate (the guidelines can easily be quantified by applying the methods described in this article to specific beam shapes).

Signal integration reduces the fluctuation-induced error from its instantaneous value with a time constant $T_{1/2} \simeq 2 \Delta_{1/2} h_v/v \sin^2 E$. Because of its dependence on beam width and elevation, the effect of signal integra-

¹⁷ R. Linfield, Tracking Systems and Applications Section, personal communication, Jet Propulsion Laboratory, Pasadena, California, 1991.

¹⁸ The measured gain of the current J and D series WVR's drift at a rate that causes an approximate 0.3-percent gain change in 1000 sec (Footnotes 4 and 5). Therefore, the stability of these WVR's should be improved by a factor of at least 60, from $\dot{g}/g \simeq 3 \times 10^{-6} \text{ sec}^{-1}$ to $\dot{g}/g \simeq 5 \times 10^{-8} \text{ sec}^{-1}$.

¹⁹ See Footnote 2.

²⁰ See Footnote 4.

tion has been neglected in Fig. 10. For example, for a WVR beamwidth $\Delta_{1/2} = 0.1$ deg (FWHM = 0.2 deg) and $E = 6$ deg, $T_{1/2}$ is 1 min, and the path delay error when $N_{v,obs} = 2$ g/cm² is about 4 mm. Hence, integrating the WVR signal over a 2-min period (which is a typical VLBI integration time) will reduce the fluctuation error to about 1 mm (which is the desired accuracy for the estimated path delays). It has been suggested that copointing a radio telescope and a WVR antenna in the same direction would simplify WVR antenna steering. The copointing will introduce a systematic error and increase the random error in the estimated path delay. For $E > 6$ deg, these additional errors will be smaller than the random error for the brightness centroid pointing for all beam sizes $\Delta_{1/2} < 1$ deg. For $\Delta_{1/2} > 1$ deg, the additional errors will increase $\propto \Delta_{1/2}^2$, and, in addition, the required WVR signal integration time to average out the fluctuation-induced error becomes longer²¹ than the $T_{1/2}$ for the brightness centroid pointing (and longer than the VLBI integration time of about 2 min). Therefore, the ability for the wide beam WVR's to be pointed at a slightly higher elevation than the radio telescope is important.

²¹ See Footnote 12.

From the brief discussion of various error sources in Section II, the biggest error in the estimated path delay is at the present time due to inaccurate modeling of the absorptivity of water vapor and uncertainties in the distribution of atmospheric parameters along the observed lines of sight (atmospheric noise). These uncertainties cause scale errors (see Fig. 1) in the estimated path delay: a systematic error of about 10 percent due to the error in the absorptivity model and between 2 and 4 percent random error due to atmospheric noise. To satisfy the 1-mm path delay accuracy requirement at a 6-deg elevation, the atmospheric noise must be reduced to the 0.17-percent level by, for example, custom tailoring the retrieval algorithm coefficients to specific sites and a set of observing conditions.

The accuracy of the present absorptivity models could be improved by better modeling and model calibration, using for example, a comparison of WVR and radiosonde data, direct measurements or estimates using interferometric data reduction of atmospheric path delays, or measurement of water vapor absorptivity in a laboratory-controlled environment. Investigation of some of these possibilities, including that of developing mathematical methods to filter out the effect of the systematic scale error during VLBI data reduction, is part of an ongoing effort.

Acknowledgments

The author wishes to thank R. N. Treuhaft and R. P. Linfield for their many valuable comments and critical reading of the manuscript. The author is also grateful to G. M. Resch for useful discussions.

References

- [1] R. N. Treuhaft and G. E. Lanyi, "The Effect of the Dynamic Wet Troposphere on the Radio Interference Measurements," *Radio Science*, vol. 22, p. 251, 1987.
- [2] R. N. Treuhaft and S. T. Lowe, "A Measurement of Planetary Relativistic Deflection," *The Astronomical Journal*, vol. 102, p. 1879, 1991.
- [3] J. W. Armstrong, "Advanced Doppler Tracking Experiments," *Proc. of Workshop, NASA Conference on Relativistic Gravitational Experiments in Space*, Publication 3046, Annapolis, Maryland, June 28-30, 1988.
- [4] G. M. Resch, "Inversion Algorithm for Water Vapor Radiometers Operating at 20.7 and 31.4 GHz," *TDA Progress Report 42-76*, vol. October-December 1982, Jet Propulsion Laboratory, Pasadena, California, pp. 12-26, February 15, 1983.
- [5] A. J. Thompson, J. M. Moran, and G.W. Swenson, *Interferometry and Synthesis in Radio Astronomy*, New York: John Wiley and Sons, 1986.
- [6] W. C. Hamilton, *Statistics in Physical Science*, New York: The Ronald Press, 1964.

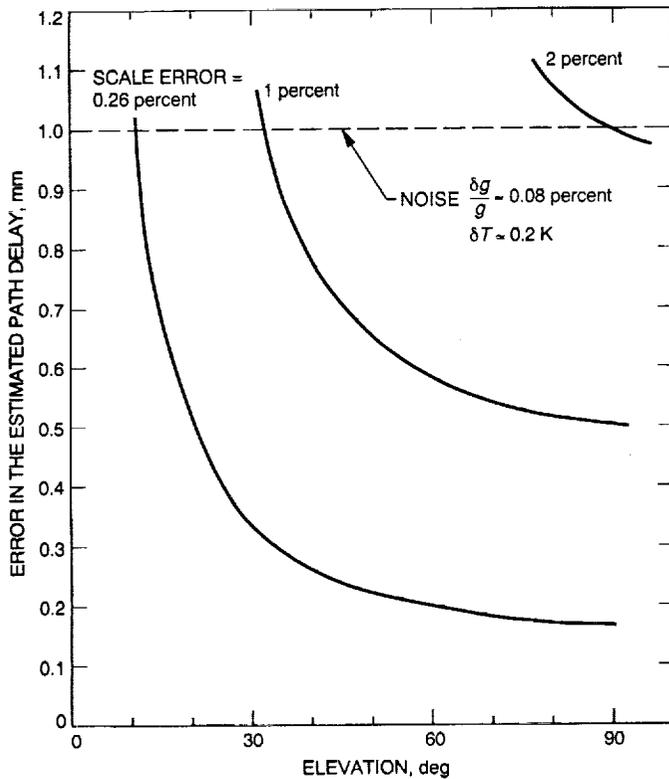


Fig. 1. The effect of stochastic (e.g., due to thermal noise and WVR gain variation) and scale (e.g., from absorption coefficient and WVR calibration) errors on path-delay estimates. The scale errors plotted are per 1 g/cm^2 of water vapor zenith column density (or, equivalently, per 6 cm of zenith path delay). The goal is 1-mm path-delay calibration accuracy in most weather conditions. (This figure assumes no reduction in the effects of a scale error from parameter estimation.)

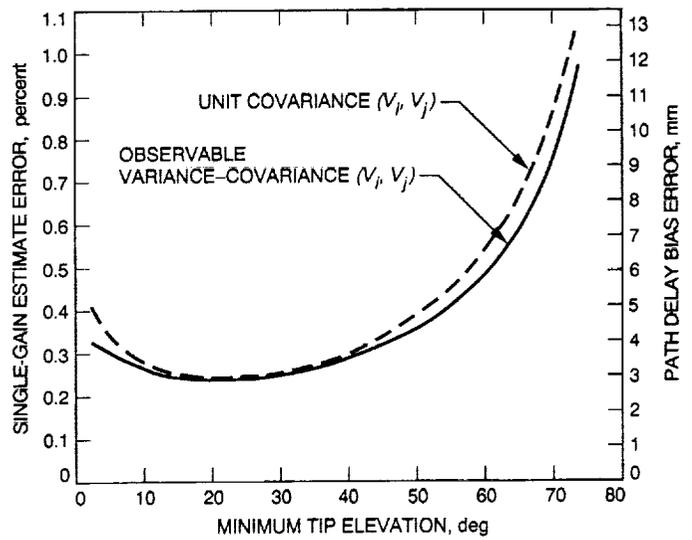


Fig. 2. The standard deviation of a single gain estimate from tip curves versus minimum tip elevation. The heavy and broken line curves were calculated using the actual and unit observable covariance-variance matrix W^{-1} in the least-squares analysis. The path delay shown is the bias error ΔL_V .

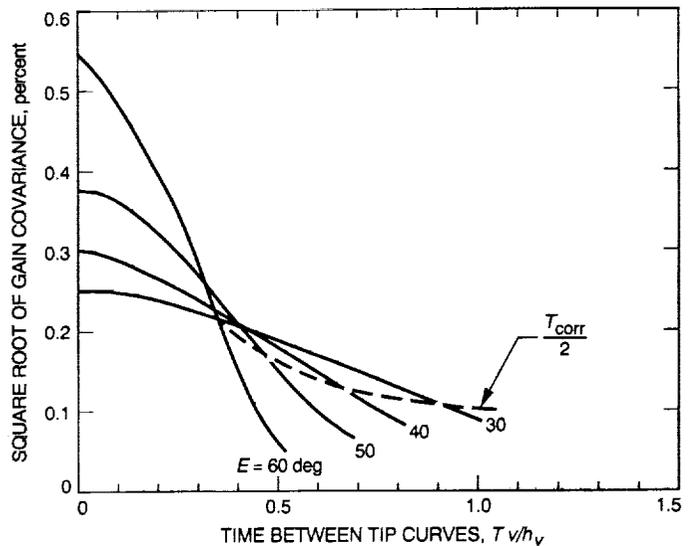


Fig. 3. Temporal development of the square root of the covariance of successive gain estimates versus time, T , elapsed between tip curves. The single gain estimates decorrelate within the time $T_{corr} = h_v/(v \tan E_{min})$. The broken line curve connects covariance values at $T = T_{corr}/2$.

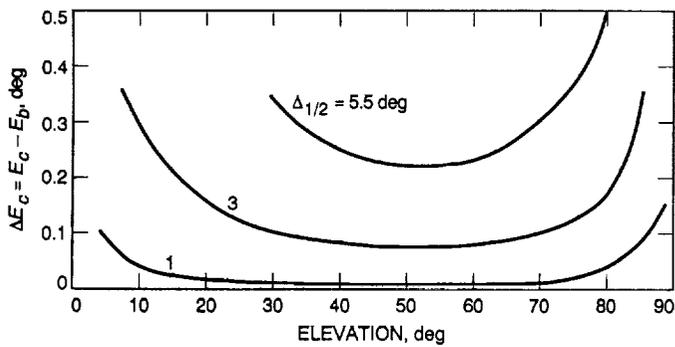


Fig. 4. The offset (ΔE_c) between elevations of the WVR beam geometrical center and brightness centroid. For small beam widths, $\Delta E_c \propto \Delta_{1/2}^2$. Note that as E_c approaches 90 deg, the offset increases to $\Delta_{1/2}/\sqrt{3}$.

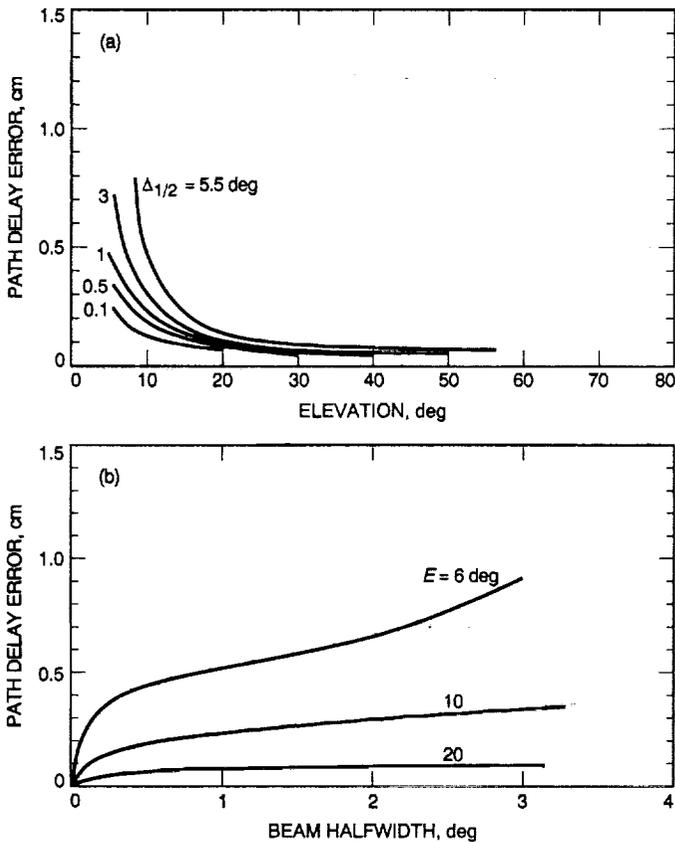


Fig. 5. Path delay stochastic error for the brightness centroid pointing due to WVR beam averaging of tropospheric fluctuations (instantaneous error, per 1 g/cm^2 of zenith water vapor column density): (a) versus elevation and (b) versus beam half-width.

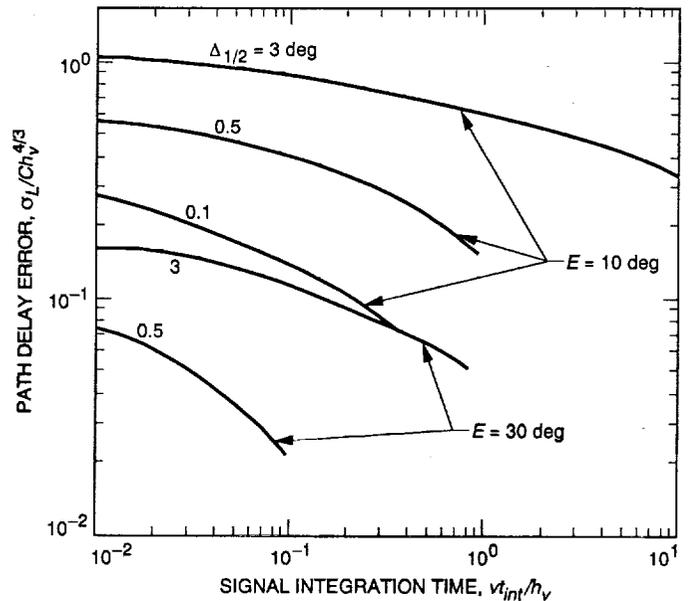


Fig. 6. Path delay stochastic error for the brightness centroid pointing versus the signal integration time. For a symmetrical radiation pattern, the error depends only very little on the wind direction (for different wind directions, the calculated errors differed by less than 15 percent).

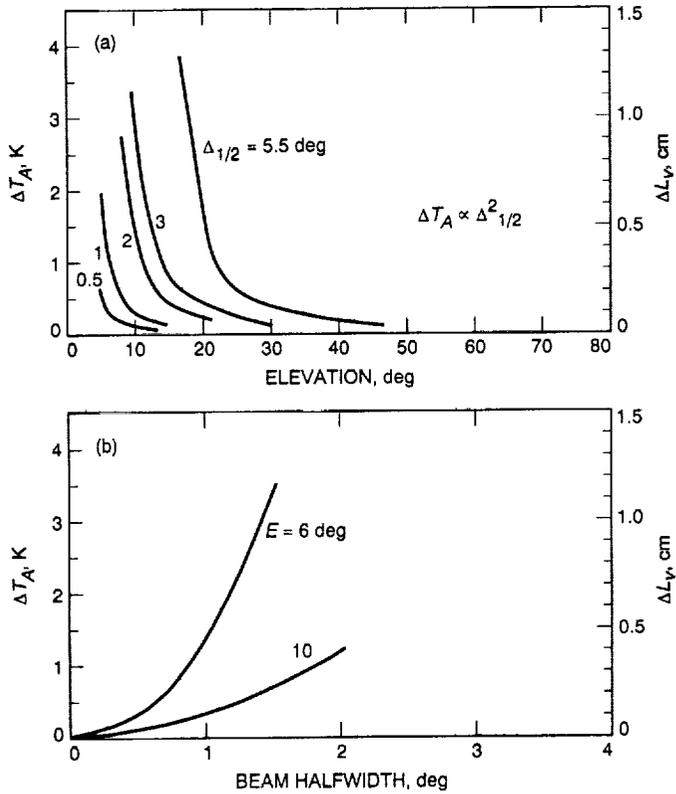


Fig. 7. The systematic difference ΔT_A between the antenna temperature and the brightness temperature at 20.6 GHz (and the corresponding path delay error) in the beam geometrical center due to WVR beam width averaging of air mass (per 1 g/cm^2 of zenith water vapor column density): (a) versus elevation and (b) versus beam half-width.

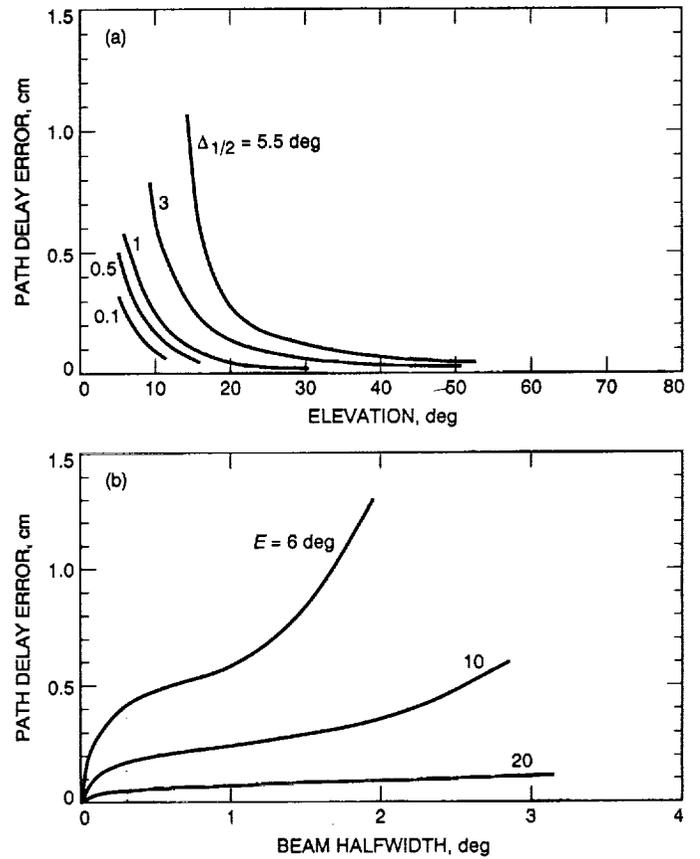


Fig. 8. Path delay stochastic error for the geometrical center pointing due to WVR beam averaging of tropospheric fluctuations (per 1 g/cm^2 of zenith water vapor column density): (a) versus elevation and (b) versus beam half-width.

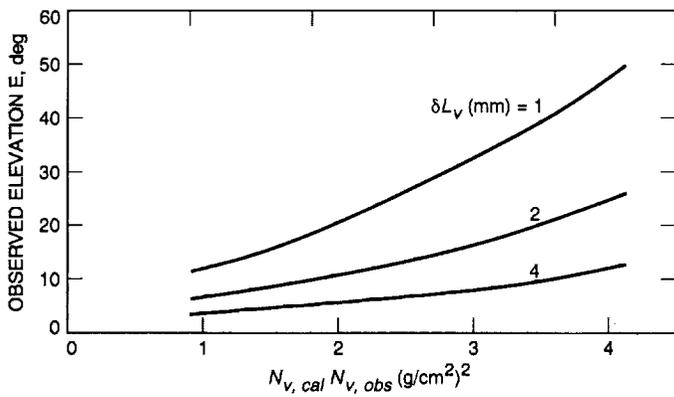


Fig. 9. Curves of path delay scale errors (δL_v) for a single gain estimate. For each curve, the error is less than the cutoff error above and to the left of the curve. The values of $N_{v,cal}$ and $N_{v,obs}$ are zenith water vapor column densities during WVR calibration and radiometric observation, respectively.

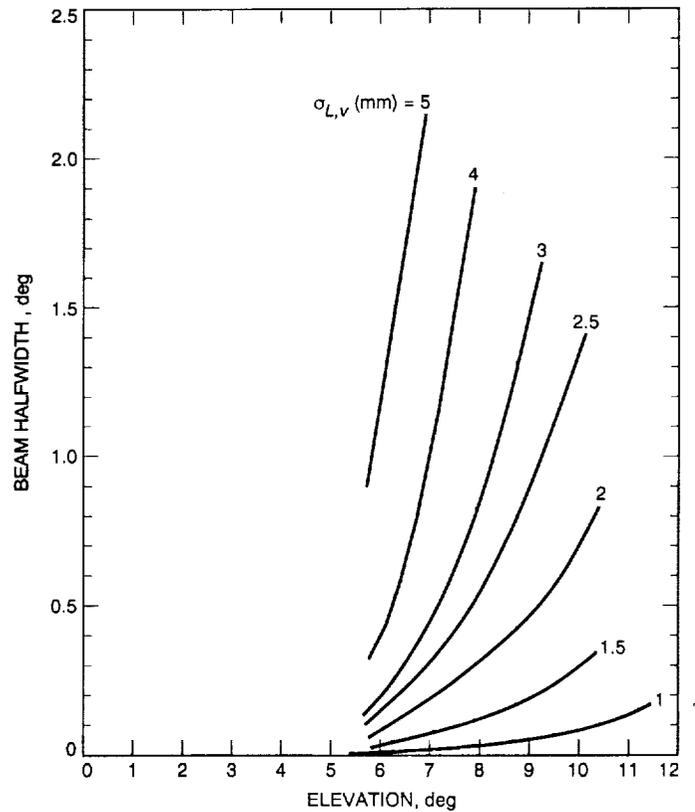


Fig. 10. Curves of tropospheric fluctuation-induced path delay error ($\sigma_{L,v}(E_b)$) for the brightness centroid (per 1 g/cm^2 of zenith water vapor column density) due to WVR beam averaging. In more humid weather, the errors scale with humidity.

Appendix A

The Method of Least Squares

For brevity of notation, the voltage recorded by a WVR pointed in a direction E_i is designated V_i , and the line-of-sight opacity is designated τ_i . The duration of a single tip sequence is less than the time it takes to decorrelate the tropospheric inhomogeneities. The i th tip curve voltage, V_i , is then modeled as

$$\begin{aligned} V_i &= g(T_C e^{-\tau_i} - T_M(1 - e^{-\tau_i}) - T_{ref}) \\ &\simeq g \tau_i T_{MC} - g T_{RC} \end{aligned} \quad (\text{A-1})$$

where $T_{RC} = T_{ref} - T_C$, $T_{MC} = T_M - T_C$, and where the linearized Eq. (A-1) is a good approximation to the full radiation transport equation for most ($\tau_i < 0.5$) tropospheres of interest.

In tip curve analyses, the WVR recorded data are fit by assuming a temporally constant and spatially homogeneous stratified troposphere. For $i = 1, \dots, N$ data, where N is the number of tip curve elevations, the V_i 's represent a set of N equations for two solve-for parameters, \hat{g} and $\hat{\tau}_z$. Mapping the τ_i 's to zenith by using $\tau_i = \hat{\tau}_z A_i$, where $\hat{\tau}_z$ is the estimate for the zenith opacity τ_z , and the air mass $A_i = 1/\sin E_i$, the equations are solved by the method of least squares [6], as follows.

Defining solve-for parameters and observable column vectors $X = [g \ \tau_z, g]$ and $V = [V_1, \dots, V_N]$, respectively, the design matrix \mathcal{A} (\mathcal{A} has dimensions $N \times 2$) in $V = \mathcal{A} X$ is

$$\mathcal{A} = \begin{pmatrix} T_{MC} A_1 & -T_{RC} \\ T_{MC} A_2 & -T_{RC} \\ \vdots & \vdots \\ T_{MC} A_N & -T_{RC} \end{pmatrix} \quad (\text{A-2})$$

Assuming that the errors in V_i 's have zero means and a variance-covariance matrix W^{-1} , minimization of the

quadratic form $((V - \mathcal{A}\hat{X})^T W (V - \mathcal{A}\hat{X}))$ yields the following estimates for g and τ_z :

$$\hat{X} = (\mathcal{A}^T W \mathcal{A})^{-1} \mathcal{A}^T W V \quad (\text{A-3})$$

where \hat{X} is the column vector $\hat{X} = [g \hat{\tau}_z, \hat{g}]$, and the superscript T designates the transpose matrix. By substituting Eq. (A-1) into Eq. (A-3), one can easily verify that the statistically averaged estimated gain is equal to the WVR gain (i.e., $\langle \hat{g} \rangle = g$), as it should be. The estimated gain standard deviation is the square root of $\sigma_{\hat{g}}^2 \equiv \text{cov}(\hat{g}, \hat{g})$, which is the matrix element $(\sigma_{\hat{X}}^2)_{2,2}$ of

$$\begin{aligned} \sigma_{\hat{X}}^2 &\equiv \text{Exp} \{ (\hat{X} - X)(\hat{X} - X)^T \} \\ &= B^{-1} \mathcal{A}^T W \text{cov}(V, V^T) W \mathcal{A} B^{-1} \end{aligned} \quad (\text{A-4})$$

where Exp designates the expectation value, $\text{cov}(V, V^T)$ is the actual observable covariance-variance matrix, and $B = \mathcal{A}^T W \mathcal{A}$. Equations (A-3) and (A-4) yield \hat{g} , $\hat{\tau}_z$, and $\sigma_{\hat{X}}^2$ for given W_{ij} and $\text{cov}(V, V^T)$.

In practice, \hat{X} and $\sigma_{\hat{X}}^2$ can be derived either by using some assumed W^{-1} (the so-called consider analysis [1]), or by setting W^{-1} equal to the observable variance-covariance matrix $\text{cov}(V, V^T)$. The most common (and simplest) form of the assumed W^{-1} is the unit matrix $(W^{-1})_{i,j} = \delta_{i,j}$, where $\delta_{i,j}$ is the Kronecker delta. Taking the unit W^{-1} corresponds to assuming that the observable errors are uncorrelated with equal variances and reduces the minimization procedure to the minimization of the sum of squares. The latter, i.e., setting $(W^{-1})_{i,j} = \text{cov}(V_i, V_j)$, minimizes the variance-covariance matrix $\sigma_{\hat{X}}^2$ as $\sigma_{\hat{X}}^2 = (\mathcal{A}^T W \mathcal{A})^{-1}$.

The results of these calculations are discussed in the main text. Evaluation of $\text{cov}(V_i, V_j)$ using the Kolmogorov turbulence model is described in Appendix B.

Appendix B

Evaluation of the Correlation Functions Using the Kolmogorov Turbulence Model

To evaluate correlations between the simulated data, the voltage V_i associated with the sky brightness $T_{B,i}$ in the elevation direction E_i is first related to the line-of-sight opacity τ_i by using Eq. (A-1). This expresses correlations between V_i 's in terms of correlations between τ_i 's. Next, neglecting fluctuations in the dry component of τ_i , correlations between wet opacities are expressed in terms of correlations of wet refractivity by expressing the wet opacity ($\tau_{v,i}$) as the line-of-sight integral

$$\tau_{v,i} = \int_0^{h_v / \sin E_i} \alpha_v(\vec{r}_i) dr_i \simeq \frac{\tau_{v,z}}{L_{v,z}} \int_0^{h_v / \sin E_i} \chi(\vec{r}_i) dr_i \quad (\text{B-1})$$

where h_v is the height of the tropospheric slab; α_v is the water vapor absorptivity per unit length; $\tau_{v,z}$ and $L_{v,z}$ are the average values of wet opacity and path delay at zenith, respectively; and $\chi(\vec{r}_i)$ is the index of refraction - 1. The correlations between $\chi(\vec{r}_i)$'s are evaluated by using the Kolmogorov turbulence structure function, Eq. (8) of the main text. In what follows, the above described procedure of evaluating the correlations is exemplified in the evaluation of the observable variance-covariance matrix W^{-1} .

By using Eq. (A-1), the matrix element $W_{i,j}^{-1}$ is written as

$$\begin{aligned} W_{i,j}^{-1} &\equiv \text{cov}(V_i, V_j) = g^2 T_{MC}^2 \text{cov}(\tau_{v,i}, \tau_{v,j}) \\ &= g^2 T_{MC}^2 (\langle \tau_{v,i}, \tau_{v,j} \rangle - \langle \tau_{v,i} \rangle \langle \tau_{v,j} \rangle) \end{aligned} \quad (\text{B-2})$$

where $\langle \dots \rangle$ signifies statistical ensemble averaging.

By using Eq. (B-1) and the expression (Eq. A.3 of [1])

$$\langle \chi_i \chi_j \rangle = \langle \chi^2 \rangle - \frac{1}{2} D_\chi(\vec{r}_i - \vec{r}_j) \quad (\text{B-3})$$

where $\chi_i \equiv \chi(\vec{r}_i)$ and $D_\chi(\vec{r}_i - \vec{r}_j) \equiv \langle (\chi_i - \chi_j)^2 \rangle$ is the refractivity spatial structure function by interchanging the

order of integration and ensemble averaging, and then setting $dr_i = A_i dz$ and $dr_j = A_j dz'$, Eq. (B-2) yields

$$\begin{aligned} W_{ij}^{-1} &= \left(\frac{\tau_{v,z} g T_{MC}}{L_{v,z}} \right)^2 A_i A_j \\ &\times \int_0^{h_v} dz \int_0^{h_v} dz' \left(\sigma_\chi^2 - \frac{D_\chi(\vec{r}_i - \vec{r}_j)}{2} \right) \end{aligned} \quad (\text{B-4})$$

where the variance σ_χ^2 of the wet refractivity fluctuation is independent of spatial coordinates and is obtained by letting the distance R go to infinity in

$$\sigma_\chi^2 \equiv (\langle \chi^2 \rangle - \langle \chi \rangle^2) = \frac{D_\chi(R = \infty)}{2} = \frac{N_v^2 C^2 L_s^{2/3}}{2} \quad (\text{B-5})$$

where the last expression on the right-hand side has been obtained by evaluating the asymptotic $D_\chi(R = \infty)$ by using Eq. (8) of the main text, and where N_v is the water vapor column density at zenith in g/cm^2 , and L_s is the tropospheric turbulence saturation length. The reason why $D_\chi(\infty)$ should converge as R becomes very large has been discussed in [1].

The covariance of successive gain estimates and the effect of signal integration on beam averaging of tropospheric fluctuations involve evaluation of correlations of type $\langle \tau_i(t) \tau_j(t+T) \rangle$. By using the "frozen" troposphere model [1], these correlations were evaluated using the expression

$$\langle \chi(\vec{r}_i, t) \chi(\vec{r}_j, t+T) \rangle = \langle \chi^2 \rangle - \frac{1}{2} D_\chi(\vec{r}_i - \vec{r}_j + \vec{v} T) \quad (\text{B-6})$$

where the structure function

$$D_\chi(\vec{r}_i - \vec{r}_j + \vec{v} T) = \langle (\chi(\vec{r}_i, t) - \chi(\vec{r}_j, t+T))^2 \rangle$$

is the same as Eq. (8) of the main text.

54-32

128437

N93-19417
Q-11

The Goldstone Real-Time Connected Element Interferometer

C. Edwards, Jr., D. Rogstad, D. Fort, L. White, and B. Iijima
Tracking Systems and Applications Section

Connected element interferometry (CEI) is a technique of observing a celestial radio source at two spatially separated antennas and then interfering the received signals to extract the relative phase of the signal at the two antennas. The high precision of the resulting phase delay data type can provide an accurate determination of the angular position of the radio source relative to the baseline vector between the two stations. This article describes a recently developed connected element interferometer on a 21-km baseline between two antennas at the Deep Space Network's Goldstone, California, tracking complex. Fiber-optic links are used to transmit the data to a common site for processing. The system incorporates a real-time correlator to process these data in real time. The architecture of the system is described, and observational data are presented to characterize the potential performance of such a system. The real-time processing capability offers potential advantages in terms of increased reliability and improved delivery of navigational data for time-critical operations. Angular accuracies of 50-100 nrad are achievable on this baseline.

I. Introduction

Interferometric techniques have been used for several decades in the astronomy community to obtain very high angular resolution images or to determine astrometric positions of celestial radio sources [1-4]. By cross-correlating signals received at two spatially separated sites, one synthesizes an effective aperture corresponding to the spatial separation of the two antennas, with a resulting improvement in angular resolution. Just as the beamwidth of a

single-aperture antenna scales as λ/D , where λ is the observing wavelength and D is the antenna diameter, the resolution of an interferometer scales as λ/L , where L is the distance between the interferometer antenna pair. As will be demonstrated, this same high resolution can be used to help track and navigate interplanetary spacecraft.

Two-way tracking of the round-trip delay and Doppler shift of radio signals between Earth and a spacecraft pro-

vides a direct measurement of the geocentric range and radial velocity of the spacecraft. For current-generation spacecraft incorporating X-band (8.4-GHz) radio links, spacecraft range can be determined to 10 m, and spacecraft radial velocity to 0.1 mm/sec.

The plane-of-sky angular coordinates of the state vector, however, are more difficult to determine. Some angular information can be deduced from the signature of the Earth's rotation on the Doppler observable. Accuracies of about 150 nrad can be obtained from an 8- to 12-hr arc of Doppler data. Near the celestial equator, the accuracy of the declination component of angular position degrades as $1/\sin(\delta)$, where δ is the spacecraft declination, due to a singularity in the Doppler partial derivative.

Interferometric techniques have been developed and are used in the Deep Space Network to improve the ability to track angular spacecraft position. Currently, the technique of very long baseline interferometry (VLBI) is used in the DSN for spacecraft tracking. As shown in Fig. 1, VLBI measures angular position by determining the delay, τ , between arrival of the signal wavefront from a radio source at the two antennas. This delay is related to the angle, θ , between the baseline and the direction to the radio source:

$$\tau = \frac{1}{c} B \cos \theta \quad (1)$$

where c is the speed of light. Because of the large spatial separation between VLBI antennas, data are recorded at each antenna along with timing references from extremely stable clocks at each site. The recorded data are then transmitted to a common site for subsequent correlation processing, in which the delay is extracted. Final observables are typically unavailable until hours or days after the observation is complete.

The current DSN VLBI system can provide angular accuracies of about 30 nrad for spacecraft tracking, based on a short, 30-min observation, which represents a significant improvement over Doppler tracking alone. The VLBI observation provides a direct, geometric determination of the spacecraft angular position, in contrast to the Doppler case in which the angular position is extracted from the signature of Earth's rotation on the Doppler observable over a long data arc. In addition, VLBI suffers no degradation in accuracy at declinations near 0 deg. Combining information from the California-Spain and California-Australia DSN baselines allows good determination of both the right ascension and declination of the spacecraft throughout the ecliptic plane.

By rearranging Eq. (1), one finds that the angular accuracy $\delta\theta$ is related to the accuracy $\delta\tau$ with which the delay is determined and the length of the baseline projected onto the plane of the sky:

$$\delta\theta = \frac{c\delta\tau}{B \sin \theta} \quad (2)$$

From this, it is seen that angular accuracy is improved by improving the delay measurement or by increasing the baseline length. This latter point has driven the development of VLBI on intercontinental baselines; the DSN baselines from California to Spain and California to Australia are roughly 8000 and 10,000 km in length, a sizable fraction of the Earth's total diameter.

For several years, the authors have been investigating the extent to which interferometry on relatively short baselines of under 100 km in length can provide medium accuracy, 50- to 100-nrad angular tracking. To achieve this degree of accuracy on such a short baseline requires improving the precision of the interferometer delay measurement. As will be shown in the next section, on these shorter baselines one can make full use of the interferometer phase observable to achieve this gain in precision. In addition, there are a number of operational benefits to performing interferometry on a short baseline that have motivated interest in investigating connected element interferometry (CEI) [5]:

- (1) By using fiber-optic links, the data from the various antennas can be brought together to a common site for real-time correlation processing, reducing the turnaround time for delivering tracking observables to a navigation filter.
- (2) Real-time processing provides a real-time monitor of the complete interferometry system. Many problems affecting VLBI observations are not detected until correlation processing; real-time correlation would help to uncover such problems during the observation in time to correct them.
- (3) A common clock can be distributed to both antennas, allowing them to be operated coherently and eliminating the need to solve for a clock rate offset between stations.
- (4) Propagation media errors are significantly reduced due to common-mode cancellation on the short baseline.
- (5) The short baseline results in longer mutual visibility periods and higher elevation angles. This eases observation scheduling and reduces the effect of propagation media errors.

II. The CEI Phase Observable

The basic observable of any interferometer is the relative phase of a received signal at two spatially separated antennas. This phase can be thought of as a measure of the interferometer delay τ in units of the observing wavelength. One can write the interferometer phase as

$$\phi + 2\pi N = \omega_{RF} \left(\frac{1}{c} B \cos \theta + \tau_{clock} + \tau_{trop} + \tau_{ion} + \tau_{inst} \right) + \phi_{LO} \quad (3)$$

where ω_{RF} is the RF observation frequency, and where the total delay τ is composed of the geometric delay given in Eq. (1), the clock offset between stations τ_{clock} , the tropospheric and ionospheric propagation media delays τ_{trop} and τ_{ion} , and any uncalibrated instrumental delay τ_{inst} . (Each of these delay terms represents the differential effect between stations.) In addition, there is an overall unknown phase offset ϕ_{LO} between the aggregate local oscillators at each station. The term $2\pi N$ represents the cycle ambiguity associated with the phase data type. The high precision of the phase data type is not usable until this cycle ambiguity has been properly resolved.

For VLBI, uncertainties in the delay model typically prevent resolving the phase ambiguity. Instead, the delay is measured directly by determining the group delay $\partial\phi/\partial\omega$. In practice, what is actually measured is the quantity $(\phi_1 - \phi_2)/(\omega_1 - \omega_2)$ for two or more nearby frequencies. This group delay observable provides an unambiguous, but much less precise, measure of the interferometer delay. The group delay is less precise by the ratio of $\omega/\Delta\omega$, where ω is the RF observing frequency, and $\Delta\omega$ is the spanned bandwidth over which the group delay is calculated. For an X-band (8.4-GHz) spacecraft downlink with VLBI tones spanning a 40-MHz bandwidth, the phase observable is thus more than two orders of magnitude more precise than the group delay. The inability to resolve the integer cycle ambiguity prevents the use of this precise phase observable on intercontinental baselines.

On short baselines, however, the a priori delay model is sufficiently accurate to allow phase ambiguity resolution. Biases associated with the clock, instrumental, and LO terms in Eq. (3) are handled by differencing the phase observable for two sequential radio source observations. This differencing also greatly attenuates the effects of the propagation media errors if the sources are angularly close. The differential phase observable for two sources, A and B , can then be written

$$\Delta\phi + 2\pi\Delta N = \omega_{RF} \left(\frac{1}{c} B [\cos \theta_A - \cos \theta_B] + \Delta\tau_{trop} + \Delta\tau_{ion} + \Delta\tau_{inst} \right) \quad (4)$$

where ΔN now represents the differential phase ambiguity. The term $\Delta\tau_{inst}$ is retained to represent any stochastic temporal instability in the CEI signal path over the time between the two observations.

III. CEI Error Sources

CEI error analysis focuses on two distinct issues: the a priori model delay accuracy required to achieve ambiguity resolution and the final a posteriori accuracy obtained from the resolved phase observable. The ability to determine the differential phase ambiguity ΔN is dependent on the a priori model uncertainties associated with the terms on the right-hand side of Eq. (4). Basically, the overall delay model must be known to much better than a wavelength of the RF observing frequency. The DSN downlink spacecraft frequencies of 2.3 GHz (S-band) and 8.4 GHz (X-band) correspond to wavelengths of 13 and 3.6 cm, respectively. Once the ambiguity is resolved, many of these same errors will limit the final astrometric accuracy of the observation. In the following sections, each of these error sources is briefly examined.

A. Baseline

For short intracomplex DSN baselines, the vector between stations is typically known to 3–5 mm or less, based on geodetic interferometry experiments. For small angular separations between the radio sources, the impact of this uncertainty on the differential phase delay is further reduced, roughly by the angular source separation in radians. Thus, for a 10-deg separation, this error is below 1 mm.

Gravity deformation, wind loading, and thermal expansion could also potentially introduce antenna distortions at the millimeter level. Here again, the differential nature of the CEI observations is key to reducing this error source. Angularly close sources will have similar gravity deformation; similarly, differencing observations over a short time scale will help to reduce the effects of wind and thermal distortions. Further error cancellation will occur if the two antennas used in the observation are of identical size and design.

B. Source Position

Here it is assumed that one of the sources is a well-known reference quasar with an a priori position uncer-

tainty of 10 nrad and that the other source, the target source whose position is to be determined (e.g., a spacecraft), has a much larger source position uncertainty $\delta\theta$. On a 21-km baseline (corresponding to the longest available DSN intracomplex baseline), the 10-nrad error of the reference source corresponds to a 0.2-mm error. To unambiguously determine the phase ambiguity, a priori knowledge of the target source position must satisfy

$$\delta\theta < \frac{1}{6} \frac{\lambda}{B \sin \theta}$$

where λ is the observing wavelength. The factor of 1/6 ensures that anything less than a 3-sigma position error will cause less than a 1/2-wavelength error, and thus not cause a cycle ambiguity error. For S-band observations, and a projected baseline length of 21 km, this corresponds to a required a priori position knowledge of about 1 μ rad for the target source.

C. Troposphere

The troposphere error corresponds to the double difference of the tropospheric path delay along the four lines of sight from the two ground stations to the two radio sources. Most of the overall tropospheric effect cancels in this double differencing; the remainder represents the spatial and temporal fluctuations in the troposphere on a scale determined by the spatial separation of the antennas, the angular separation of the sources, the scale height of the troposphere, and the time separation of the two scans. These effects have been studied thoroughly elsewhere [6,7]; here the authors will just characterize the expected magnitude of this error source: For an angular source separation of 10 deg, a 21-km baseline, an average elevation angle of 45 deg, and a time separation of 200 sec between scans, the differential troposphere error is expected to be about 5 mm. This error grows to 10 mm when the mean elevation drops to 20 deg. Because this error is due primarily to small-scale fluctuations in the troposphere, the error is largely uncorrelated from one differential observation to the next and thus can be reduced by repeated observations.

D. Ionosphere

The ionosphere causes a dispersive phase error for each ray path of the form

$$\delta\phi_{[cyc]} = -1.34 \frac{TEC_{[10^{16} \text{el}/\text{m}^2]}}{\nu_{[\text{GHz}]}}$$

where $\delta\phi_{[cyc]}$ is the phase in cycles, $TEC_{[10^{16} \text{el}/\text{m}^2]}$ is the integrated line-of-sight total electron content, and $\nu_{[\text{GHz}]}$ is the observing frequency in gigahertz. Note the minus sign: The ionosphere actually causes the phase of the wavefront to advance. Taking the derivative with respect to frequency yields a positive group delay, as required by causality. Typical daytime values of TEC can range up to $100 \times 10^{16} \text{el}/\text{m}^2$ or more at zenith and three times higher at low elevations, which corresponds to tens of cycles at X-band and 100 cycles or more at S-band. As in the case of the troposphere, however, most of this error cancels in the double-differenced CEI observable, with the residual error being due to small-scale ionospheric inhomogeneities. While theoretical understanding of these fluctuations is limited, empirical data suggest that the ionosphere error for differential CEI observations on a 21-km baseline should be at or below the millimeter level at X-band,¹ representing just a few percent of a cycle of X-band phase. At S-band, this error grows to roughly one-tenth of a cycle. This is large enough to be of concern, but should not prevent accurate phase ambiguity resolution. Given the size of this error source and the variable nature of the ionosphere, more data on ionosphere fluctuations would be welcome. Experience gained in operating the Goldstone CEI will help to evaluate the magnitude of this error. If dual-frequency S- and X-band data are available, and no cycle errors are made at either band, this dispersive error source is eliminated by forming the appropriate linear combination of S- and X-band phase delay observables, which eliminates the charged particle error:

$$\Delta\tau_{S/X} = \left(\frac{\omega_X^2}{\omega_X^2 - \omega_S^2} \right) \frac{\Delta\phi_X}{\omega_X} - \left(\frac{\omega_S^2}{\omega_X^2 - \omega_S^2} \right) \frac{\Delta\phi_S}{\omega_S}$$

E. Clocks and Instrumentation

Because a single clock is used for both stations in CEI, there is no clock rate error as in VLBI. However, there is still a clock epoch uncertainty, since the propagation delays through the frequency distribution system and through the CEI signal path itself are not calibrated at the level of an RF wavelength. Thus, the CEI phase observable contains a phase bias, which corresponds to the unknown relative phase of the local oscillators at the two stations. This bias is removed by forming differencing phase observables for two radio sources. This differencing also serves to eliminate or reduce any other biaslike errors.

¹ A. J. Mannucci, "Temporal Statistics of the Ionosphere," JPL Interoffice Memorandum 335.1-90-056 (internal document), Jet Propulsion Laboratory, Pasadena, California, October 25, 1990.

While there is no explicit clock rate term, there are instabilities in the frequency distribution system, which can lead to an apparent difference in the instantaneous reference frequency at the two stations. The delay error induced by such an instability is $T\sigma_y(T)$ where T is the time between two scans used to form a differential observable and $\sigma_y(T)$ is the Allan standard deviation on that time scale. For $T = 300$ sec and $\sigma_y(T) = 10^{-14}$, this leads to a 3-psec error; the fiber-optic clock transfer between stations is thought to be even better than this [8].

IV. Non-Real-Time Phase Delay Observations

A number of experiments have been performed on baselines within the Goldstone complex in a non-real-time mode, with data recorded separately at each station [5,9]. The goals of these experiments were to demonstrate reliable phase ambiguity resolution and evaluate the potential angular accuracy of CEI observations. The results are briefly reviewed here.

Figure 2 shows the location of existing antennas within the Goldstone, California, Deep Space Communications Complex (DSCC). Baselines of up to 21 km are available. Fiber-optic cables have been installed linking the various antenna pairs [8]. For these non-real-time experiments, such fibers were used to distribute a common frequency reference to each antenna, which allowed the separate stations to be operated coherently.

Data have been collected on the 6-km DSS-12-DSS-13 baseline and on the 21-km DSS-13-DSS-15 baseline; only the latter results are discussed here. Data were recorded at each antenna for these non-real-time experiments; correlation processing and postprocessing were performed subsequently at JPL. To simulate differential quasar-spacecraft observations, pairs of quasars were observed with angular separations of up to 20 deg. The dual frequency S-/X-calibrated observations were then used to determine the relative angular position of each quasar pair.

Ambiguity resolution was carried out as follows:

- (1) An a priori delay model was calculated for each quasar observation; an ambiguous residual phase was then calculated for the S- and X-band observation relative to this model.
- (2) The residual phases were differenced between adjacent scans for a given quasar pair. This serves to eliminate the unknown phase bias between stations and reduce many other errors through common-mode cancellation.

- (3) The a priori delay model was used to resolve the S-band phase ambiguity for the differential phase observable. In effect, this required the a priori model for the differential quasar delay to be good to about 6.5 cm or better (1/2 the S-band wavelength).
- (4) The S-band residual was then used to resolve the X-band cycle ambiguity. (The X-band ambiguity was chosen so that the S- and X-band phase delay residuals agreed to within half an X-band cycle.) This approach assumes that the dominant errors are nondispersive, which implies that the S- and X-band phase delay residuals should be the same. This permits successful X-band ambiguity resolution even when nondispersive errors are more than half an X-band cycle.

Statistical analysis of the resulting phase residuals supports the reliability of the S- and X-band phase ambiguity resolution: The raw phases cluster about integer values of the cycle ambiguity, and the width of the distribution is much less than half a cycle.

After ambiguity resolution, the S-/X-calibrated phase delay observable is formed to remove the effects of charged particles. To assess the accuracy of these phase delay observations, the data were fit to estimate the angular position of one of the quasars relative to the other. Figure 3 shows the resulting adjustment in right ascension and declination for the radio source CTA 102 for roughly 3 hours of data. The data were weighted based on a model of tropospheric fluctuations above a 21-km baseline; the observed phase delay residuals were consistent with these data weights. The semi-minor axis of the source position error ellipse is 73 nrad. While each individual measurement only provides information for one component of the plane-of-the-sky position, the baseline rotation over the full observation period allowed some determination of the orthogonal component. A second orthogonal baseline would allow better determination of both components of sky position in a short observation period.

V. Development of Real-Time Capability

Having demonstrated the capability of resolving the carrier phase ambiguity and obtaining 50- to 100-nrad angular accuracies, the subsequent goal is to demonstrate the capability to collect and process CEI data in real time. The two key components required to achieve this goal are a communications channel to bring the observed data to a common site and a real-time correlator to process the received signals. Over the past year, the authors have im-

plemented these components at Goldstone. Figure 4 shows a block diagram of the entire Goldstone real-time CEI system. The existing Mark III wide-channel-bandwidth VLBI data acquisition terminals (DATs) at each station are used to perform the downconversion, sampling, 1-bit quantization, time tagging, and formatting of the radio signal at each site. However, instead of recording the resulting bit stream on tape, the signals from DSS 15 are sent via a digital fiber-optic link to DSS 13. There, a real-time correlator, based on the architecture of the non-real-time JPL/Caltech Block II VLBI correlator, receives the signals from both stations and performs cross-correlation processing.

A. Fiber-Optic Data Link

The fiber-optic link, shown in Fig. 5, consists of two main units: a transmitter unit and a receiver unit. Commercially available equipment was used to the fullest extent possible to minimize development costs. The optical fiber itself exhibits very low dispersion and low losses. The measured end-to-end attenuation of the optical signal from DSS 15 to DSS 13 is -16 dB. Previous dispersion measurements of a 14-km length of this fiber placed an upper limit of 1 deg of phase nonlinearity over a 50-MHz bandwidth [8].

The transmitter unit is located at the remote site where it collects the Mark III data and transmits them through the fiber-optic cable to the site where the real-time correlator is located. The transmitter is composed of three main building blocks. The first of these, the interface circuitry, converts from the balanced emitter-coupled logic (ECL) voltage levels of the Mark III DAT to transistor-transistor logic (TTL) levels. The second stage of processing, based on the AMD TAXI AM7968 high-speed multiplexor chip, takes the parallel VLBI data (14 channels \times 4.5 Mbit/sec) and encodes them by using a 4/5 encoding scheme, converting them to a serial data stream with synchronization words inserted. Finally, this serial bit stream is sent to a laser module where the signal is converted to light levels sent across the optical fiber to the receiver site. The optical transmitter (PCO DTX-13-565) modulates the data onto a 1300-nm optical carrier signal, with -3 -dBm power. The aggregate bit rate on the fiber-optic link is 125 Mbit/sec.

The receiver is also composed of three main building blocks. The first of these is the pin diode receiver/comparator (PCO RTX-13-565) with -33 -dBm sensitivity; here the signal is converted back to an ECL signal. This signal is then sent to a serial-to-parallel converter/decoder (AMD AM7968) where the bit stream is

synchronized, decoded (5/4), and made available as a parallel word, along with strobe and status bits. The last stage converts the TTL signals back into the balanced ECL signals required by the correlator.

B. Real-Time Correlator

The real-time correlator, dubbed Real-Time Block 2 (RTB2) is a subset of the JPL/CIT Block II VLBI processor used for non-real-time processing of VLBI data. The RTB2 provides processing for 2 stations, 1 baseline, and 14 channels, while the full JPL/CIT Block II handles up to 4 stations, 6 baselines, and 28 channels. The large wirewrap boards comprising the system are identical to those in the Block II, and the VAX software is common to both processors. The output data files are identical to those of the Block II, and all the Block II postprocessing software can be used for RTB2. The Block II itself can be arranged to be the same as RTB2 with a quick cable change, facilitating the testing of new software on the Caltech campus, rather than in the operational machine. The system is built in a single rack and, of course, has no tape playback units. A brief description of the system follows.

An overview of the correlator is shown in Fig. 6. Standard Mark III formatted data enter RTB2 on two ribbon cables, which are the same as those that would normally go to a Honeywell 9600 tape recorder in a Mark III VLBI DAT. One of these will normally come from the local formatter and the other from a remote formatter via a fiber-optic link. Data from the two sets of 14 "tracks" first enter bit synchronizers that recover the data and clock signals and then pass to a 28-by-28 crossbar switch that can be set by the user to connect any track to any correlator channel. The 14 tracks from the local station are connected to channels 1-14 and those from the remote station to channels 15-28. The data are then passed through digital delay lines that are driven by 28 separate delay models sent from the VAX to the station processor. The output of the delay lines is fed to both the tone extractor board and the cross-correlator board. On the tone extractor board, there is one tone extractor for each channel, but it is time-multiplexed to allow four different tones to be extracted from each channel, and hence there are 112 different phase polynomial tone models sent by the VAX to the station processor. The connections from the delay lines to the cross-correlator board are arranged to correlate the first 14 channels with the second 14 channels. The user can choose to correlate 14 channels with 8 lags each, 7 channels with 16 lags each, or 1 channel of 112 lags. The last case would normally be used for searching clock delay with a delay range of 28 μ sec. The cross-correlation board is driven by 14 phase models and 14 fractional delay mod-

els for each of the two antennas sent by the VAX to the correlator processor.

The RTB2 retains all the features of the Block II, including spectral domain fractional bit correction. Normal integration times are integral seconds. To the user, the system appears as a two-station 14-channel correlator. The correlator control file (or keyboard commands) are the same as would be used for configuring and controlling the Block II by using two tape transports. The output file also appears to be the result of a Block II correlation using two tape transports. The monitor/display system runs on a VAX workstation and uses a second copy of the output file sent over the EtherNet, just as in the Block II system. The display shows, in real time, the state and quality of the data being received from each station, plots of the correlation fringes for all channels versus time, plots of the phases of all tones for one channel versus time, and plots of the integrated delay and delay rate patterns for one channel. The information to be displayed is chosen by the user with a mouse-driven graphics interface. Additional real-time features, including spacecraft tone acquisition and tracking, could be added in the future.

VI. Results and Future Plans

On June 18, 1991, the first end-to-end test of the Goldstone CEI system was performed. Observations of the Magellan spacecraft were scheduled on the DSS-13-DSS-15 baseline, concurrent with a regularly scheduled Magellan telemetry pass at DSS 15. Fringes were successfully detected that corresponded to the cross-correlation of the Magellan carriers at 2.3 and 8.4 GHz. Figure 7 shows the first detected fringes.

This initial test served to verify the end-to-end processing of the Goldstone CEI system and demonstrated the digital fiber-optic link and the real-time correlation processing capabilities. To verify the angular accuracy of the observables, differential quasar pair observations were scheduled on October 3, 1991. From 03:29:00-04:17:00 GMT, the two quasars 3C 454.3 and CTA 102 were observed at S-band and X-band. These two sources are separated by 6.8 deg on the plane of the sky. Twelve 3-minute scans were scheduled, alternating between sources. Fringes were visible in real time at the RTB2 correlator during data acquisition. Postprocessing of the correlated phase data yielded the ambiguity-resolved S-/X-band calibrated phase delay residuals shown in Fig. 8. An uncertainty of 7 psec was assigned to each scan to reflect the expected level of troposphere fluctuations, based on the model of [5,6].

These residuals were then fit to estimate the relative angular position of 3C 454.3 relative to CTA 102. Given the limited duration of this data arc, the baseline projection on the plane of the sky did not rotate through a large angle, and as a result these data alone were not adequate to estimate both components of the angular separation of the sources. Instead, due to the predominantly north-south orientation of the baseline, just the declination of the source 3C 454.3 was estimated. The position of each source is known from regular VLBI observations to an accuracy of 5 nrad, providing a truth model against which the CEI determination can be compared. The twelve observations were grouped into four sets of three scans. For each A-B-A sequence, a clock epoch and rate were fit. This served to interpolate the A observations to the epoch of the B observation and removed any constant or linear error in the phase delay. For the entire set of observations, the declination of 3C 454.3 was fit, with no a priori constraint. In addition, the relative zenith troposphere and the baseline vector components were estimated, all with 1-cm a priori constraints. These parameters were included so that the final formal error in declination would reflect potential uncertainties in station-differenced troposphere and station location. The resulting adjustment to the a priori declination of 3C 454.3 was 90 ± 88 nrad. The adjustments to the troposphere and baseline vector were small, as compared with their a priori 1-cm constraints.

Further observations of close quasar pairs to demonstrate the astrometric accuracy of the CEI system will be conducted in the near future. Other potential demonstration opportunities include observing the Galileo, Magellan, or Ulysses spacecraft. In particular, accurate angular tracking of Ulysses during its Jupiter Gravity Assist, conducted in February 1992, would help to improve the Jupiter ephemeris in the radio reference frame. This, in turn, could benefit the approach navigation for Galileo, which arrives at Jupiter in 1995. To this end, CEI observations were collected at Goldstone during the Ulysses flyby of Jupiter. Analysis of these data is currently under way. Preliminary indications suggest that angular accuracies on the order of 50 nrad will be achieved.

VII. Conclusions

A real-time CEI capability has been developed and demonstrated at the Goldstone Deep Space Communications Complex on the 21-km baseline between DSS 13 and DSS 15. The key technology developments that enabled this demonstration are a high-rate digital fiber-optic link and a real-time correlation processor. The fiber-optic link

carries digitized, time-tagged data at an aggregate bit rate of 125 Mbit/sec from SPC 10 to DSS 13. The correlator, based on the JPL/Caltech Block II VLBI processor, supports cross-correlation of up to 14 2-MHz channels of Mark-III formatted VLBI data, and allows extraction of calibration and/or spacecraft tone signals. Real-time data

have been successfully acquired for both spacecraft and quasar observations. The RTB2 system displays interferometric fringe data in real time, providing verification of successful data acquisition. Differential quasar pair observations have been performed, achieving angular accuracies of under 100 nrad for less than 1 hour of data.

Acknowledgments

The authors thank George Lutes for his assistance and expertise in the area of fiber optics; Lyle Skjerve, Ben Johnson, and the entire staff of DSS 13 for their efforts in installation and data acquisition at Goldstone; and George Resch for his early encouragement in pursuing this work.

References

- [1] W. M. Baars, J. F. van der Brugge, J. L. Casse, J. P. Hamaker, L. H. Sondaar, J. J. Visser, and K. J. Wellington, "The Synthesis Radio Telescope at Westerbork," *Proc. IEEE*, vol. 61, pp. 1258-1266, 1973.
- [2] G. Davies, B. Anderson, and I. Morison, "The Jodrell Bank Radio-Linked Interferometer Network," *Nature*, vol. 288, pp. 64-66, 1980.
- [3] A. R. Thompson, B. G. Clark, C. M. Wade, and P. J. Napier, "The Very Large Array," *Astrophys. J. Suppl.*, vol. 44, pp. 151-167, 1980.
- [4] M. J. Batty, D. L. Jauncey, P. T. Rayner, and S. Gulkis, "Tidbinbilla Two-Element Interferometer," *Astron. J.*, vol. 87, p. 938, 1982.
- [5] C. D. Edwards, "Angular Navigation on Short Baselines Using Phase Delay Interferometry," *IEEE Transactions on Instrumentation and Measurement*, vol. 38, p. 665, 1989.
- [6] R. N. Treuhaft and G. E. Lanyi, "The Effect of the Dynamic Wet Troposphere on Radio Interferometric Measurements," *Radio Science*, vol. 22, pp. 251-265, 1987.
- [7] C. D. Edwards, "The Effect of Spatial and Temporal Wet-Troposphere Fluctuations on Connected Element Interferometry," *TDA Progress Report 42-97*, vol. January-March 1989, Jet Propulsion Laboratory, Pasadena, California, pp. 47-57, May 15, 1987.
- [8] G. Lutes and A. Kirk, "Reference Frequency Transmission Over Optical Fibers," *TDA Progress Report 42-87*, vol. July-September 1986, Jet Propulsion Laboratory, Pasadena, California, pp. 1-9, September 15, 1987.
- [9] C. D. Edwards, "Development of Real-Time Connected Element Interferometry at the Goldstone Deep Space Communications Complex," paper AIAA 90-2903 presented at the AIAA/AAS Astrodynamics Conference, Portland, Oregon, August 20-22, 1990.

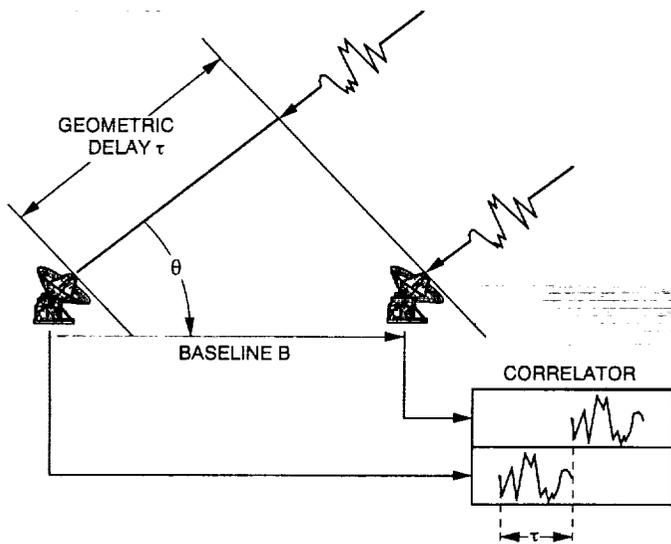


Fig. 1. Cross-correlation of signals received at two stations allows determination of the delay in arrival times of the signal wavefront at the two sites, which in turn determines the angle of the radio source relative to the baseline vector between antennas.

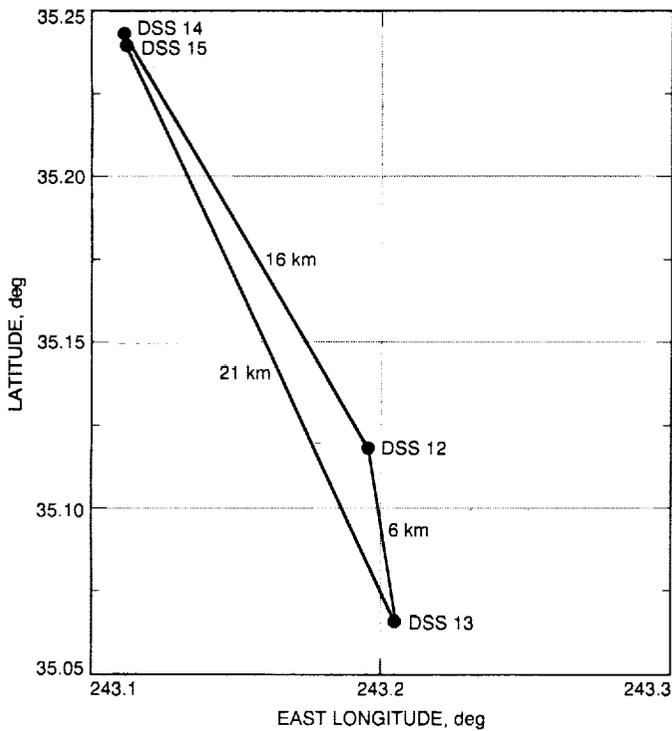


Fig. 2. Existing antennas and CEI baselines at the Goldstone DSCC.

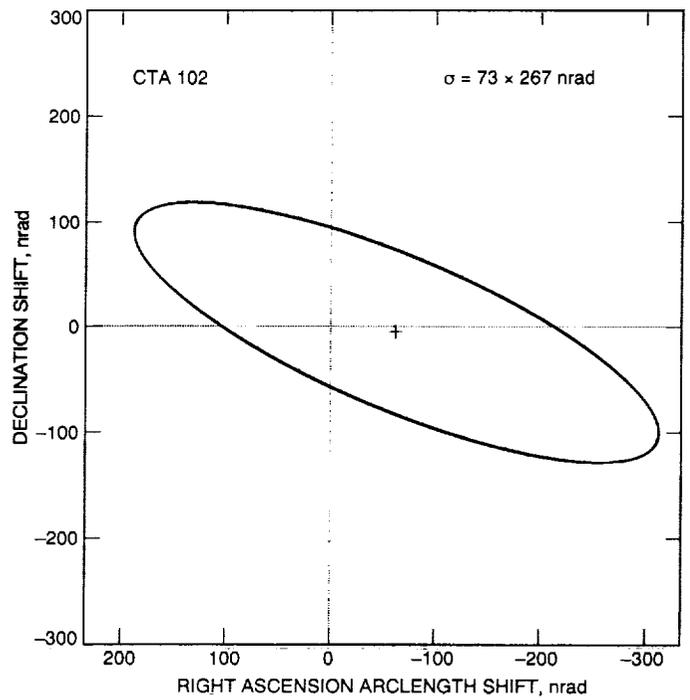


Fig. 3. Angular accuracy achieved for a non-real-time phase delay experiment on the 21-km DSS-13-DSS-15 baseline.

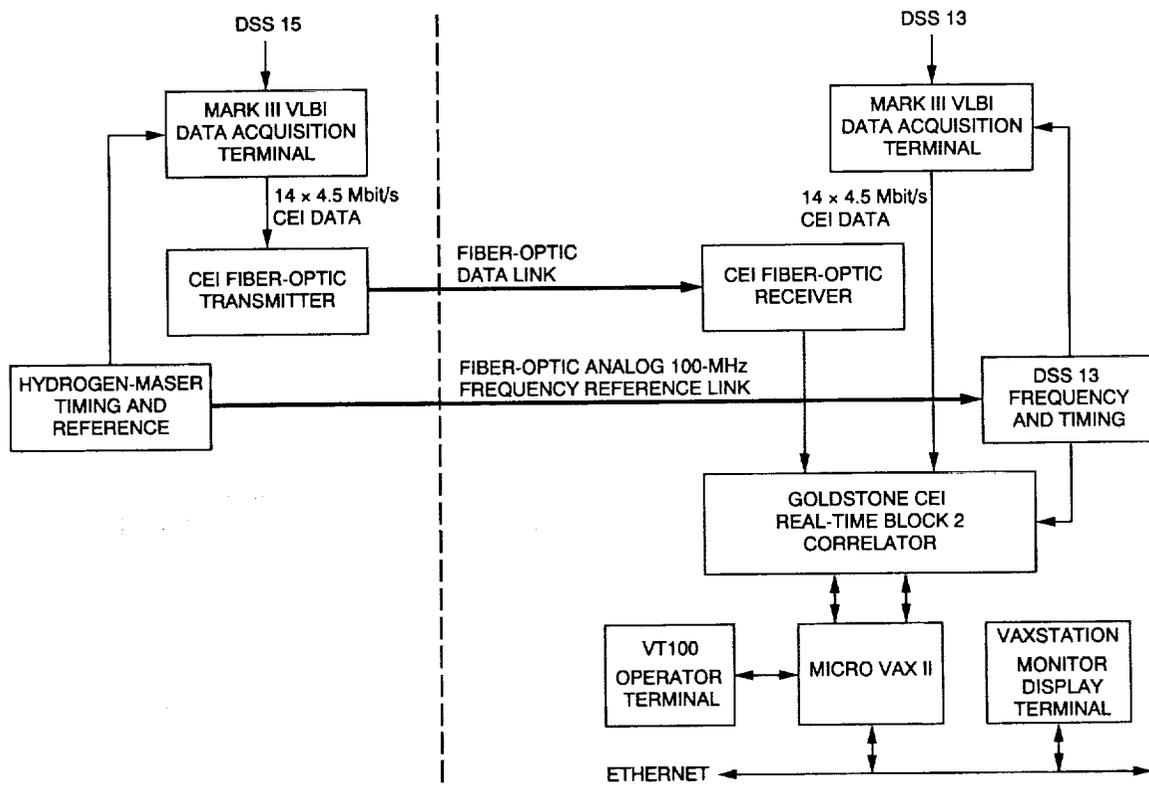


Fig. 4. Goldstone real-time CEI block diagram.

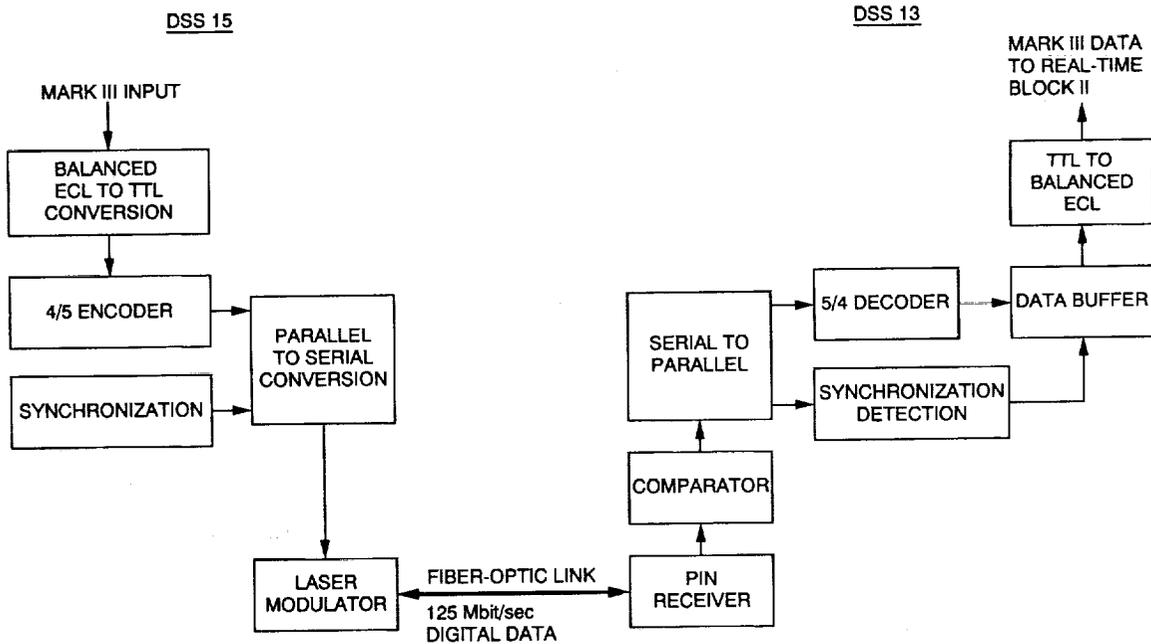


Fig. 5. Fiber-optic digital data link block diagram.

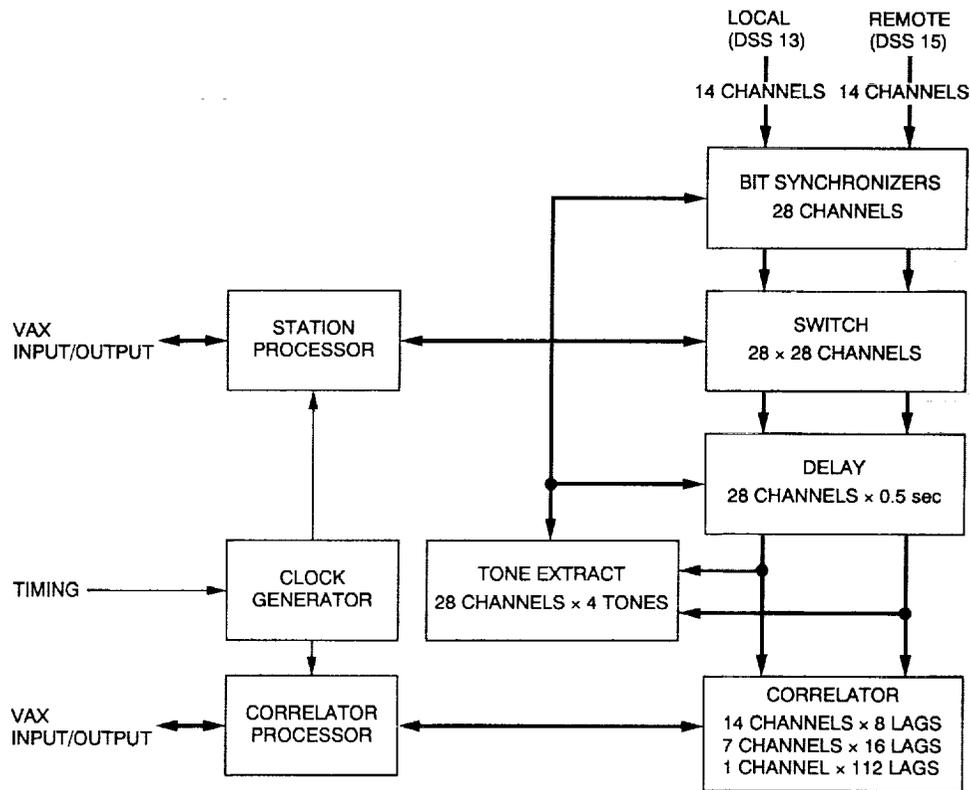


Fig. 6. Block diagram of the real-time Block II correlator.

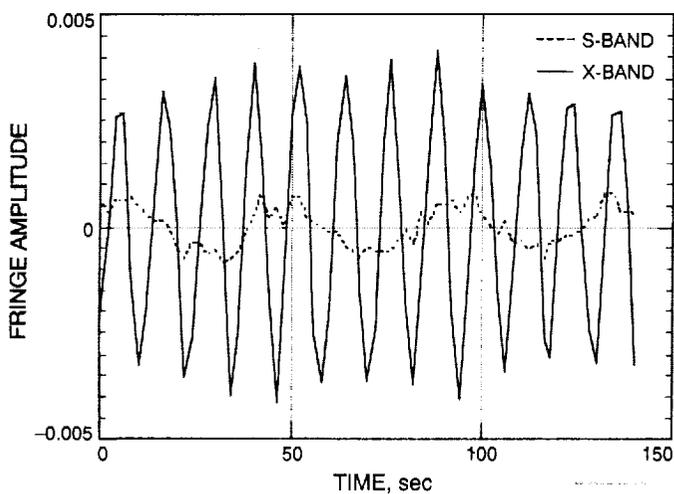


Fig. 7. First real-time fringes obtained on the Goldstone real-time CEI.

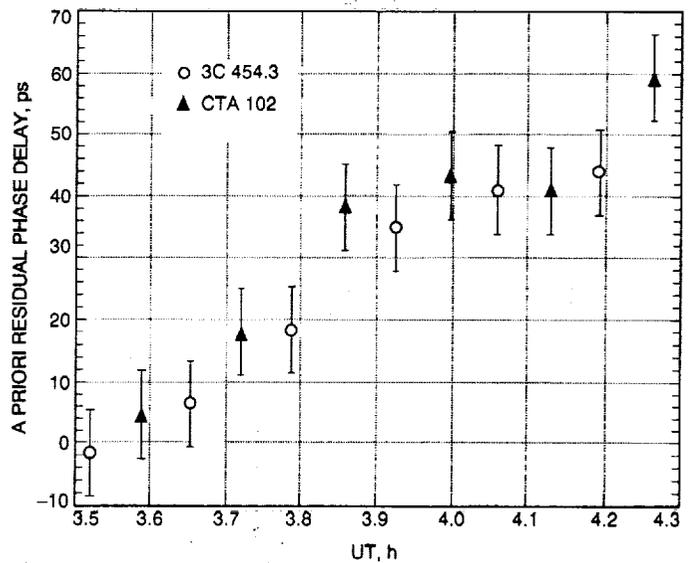


Fig. 8. Phase delay residuals for October 3, 1991, CEI observations of the radio source pair 3C 454.3-CTA 102 between DSS 13 and DSS 15.

N 93-19418

128438

p-14

Spacecraft-Spacecraft Very Long Baseline Interferometry

Part I: Error Modeling and Observable Accuracy

C. D. Edwards, Jr., and J. S. Border
Tracking Systems and Applications Section

In Part I of this two-part article, an error budget is presented for Earth-based delta differential one-way range (Δ DOR) measurements between two spacecraft. Such observations, made between a planetary orbiter (or lander) and another spacecraft approaching that planet, would provide a powerful target-relative angular tracking data type for approach navigation. Accuracies of better than 5 nrad should be possible for a pair of spacecraft with 8.4-GHz downlinks, incorporating 40-MHz DOR tone spacings, while accuracies approaching 1 nrad will be possible if the spacecraft incorporate 32-GHz downlinks with DOR tone spacings on the order of 250 MHz; these accuracies will be available for the last few weeks or months of planetary approach for typical Earth-Mars trajectories.

Operational advantages of this data type are discussed, and ground system requirements needed to enable spacecraft-spacecraft Δ DOR observations are outlined. This tracking technique could be demonstrated during the final approach phase of the Mars '94 mission, using Mars Observer as the in-orbit reference spacecraft, if the Russian spacecraft includes an 8.4-GHz downlink incorporating DOR tones. Part II of this article will present an analysis of predicted targeting accuracy for this scenario.

I. Introduction

Conventional differential very long baseline interferometry (Δ VLBI), as depicted in Fig. 1, provides angular tracking of an interplanetary spacecraft relative to one or more extragalactic radio sources (e.g., quasars). With respect to this quasar reference frame, which defines an inertial navigation reference system, Galileo-era delta differential one-way range (Δ DOR) observations between a spacecraft and an angularly nearby quasar can provide

roughly 30-nrad angular accuracy,¹ while enhancements in recorded and spanned bandwidths may enable nanoradian-level accuracy in the future [1]. However, to take full advantage of the high accuracy of Δ VLBI measurements

¹ J. Border, "Analysis of Δ DOR and Δ DOD Measurement Errors for Mars Observer Using the DSN Narrow Channel Bandwidth VLBI System," JPL Interoffice Memorandum 335.1-90-026 (internal document), Jet Propulsion Laboratory, Pasadena, California, May 15, 1990.

for planet-relative targeting, one must also have comparably accurate knowledge of the planetary ephemeris in the radio reference frame. This knowledge is currently limited to about 50 nrad for the inner planets, due mostly to uncertainty in the overall orientation of the planetary ephemerides in the radio frame, with larger ephemeris uncertainties for the outer planets. A number of techniques (observations of the millisecond pulsar PSR 1937+214, timing measurements of planetary occultations of quasars, Phobos and Magellan Δ VLBI, and intercomparison of VLBI and lunar laser ranging observations) promise to improve this knowledge to about the 25-nrad level in the next few years. Nevertheless, this still represents a large error for angular navigation, as compared with the precision of the Δ VLBI observable.

An alternative for planetary approach navigation, depicted in Fig. 2, is to use a radio signal from an orbiter or lander already at the target planet as a VLBI navigation reference beacon for the approaching spacecraft. The Δ DOR observations can then be made between the two spacecraft directly, effectively replacing the quasar with the reference spacecraft. The sequence of missions to Mars embodied in the framework of the Space Exploration Initiative (SEI) will enable such tracking opportunities, providing a number of advantages over conventional spacecraft-quasar Δ DOR. First, the frame-tie problem is eliminated: The differential measurement between the approach spacecraft and the reference spacecraft (planetary orbiter or lander) provides a direct target-relative measurement of the approach spacecraft's position. Of course, the frame-tie and ephemeris errors are replaced by any uncertainty in the reference spacecraft's position, but these errors should be much smaller than typical frame-tie or planetary ephemeris errors. For example, Doppler data typically provide subkilometer planet-relative positions for planetary orbiters, while position errors for fixed landers should be at or below about 10 meters using new differential tracking data types [2].

A second advantage of this technique is that as the approach spacecraft nears the target planet, the angular separation between the approach spacecraft and the reference spacecraft will continually decrease. This will reduce the size of a number of error sources, including platform errors (i.e., station location and Earth orientation) and propagation media effects. Thus, the highest quality data will be obtained just when it is most needed—immediately prior to orbit insertion.

Finally, because the spacecraft signals are deterministic, one-way spacecraft range observables can be generated locally at each station, without the need for the

wideband data transfer and cross-correlation between stations that is required for generating a quasar group delay observable. Only simple carrier phase tracking of several sinusoidal tones from each spacecraft is required. This has several advantages in terms of efficiency and reliability, including real-time validation of successful signal reception and real-time generation of spacecraft phase and delay observables at each station, as well as near-real-time generation of station-differenced delays. Only a very small amount of data must be brought together to form the station-differenced observables (for example, the time-tagged spacecraft phases at a one-hertz rate) and so, in principle, the spacecraft-spacecraft Δ DOR observable could be available within minutes of the actual observation for input into a navigation filter.

A goal of this article is to examine the error budget for spacecraft-to-spacecraft Δ DOR. Each physical error source is examined and, to the extent possible, parameterized as a function of angular separation. The total angular tracking error as a function of spacecraft angular separation is then calculated for spacecraft with DOR tones at either X-band (8.4 GHz) or Ka-band (32 GHz.) Possible applications of this technique are discussed. In particular, an early opportunity to demonstrate spacecraft-spacecraft Δ DOR will present itself in September 1995, the tentative arrival date of the Russian Mars '94 mission at Mars, where the U.S. Mars Observer spacecraft will already have been in orbit for over two years. To evaluate the potential navigation improvement for this Mars '94 approach scenario, a covariance analysis is performed in Section II of this article, which is based on the spacecraft-spacecraft Δ DOR error budget.

II. Error Analysis

A. Observation Description

The nominal observation scenario considered here consists of three scans: a 60-sec observation of spacecraft A, then a 120-sec observation of spacecraft B, and finally another 60-sec observation of spacecraft A, with slew times of 60 sec allowed between scans. For each observation, two stations spanning an intercontinental baseline simultaneously observe a number of sinusoidal tones (referred to as DOR tones) from one of the spacecraft, which provides a group delay measurement of the difference in arrival times of that spacecraft's signal at two tracking stations. A differential observable is then formed by interpolating the two observations of spacecraft A to the epoch of the observation of spacecraft B, thereby eliminating any errors that are linear in time, and then differencing the interpolated delay for spacecraft A from the observed delay

for spacecraft B. These observation times are long enough to provide sufficient signal-to-noise ratio (SNR) for typical DOR tone amplitudes, yet short enough to keep various stochastic errors small. Two spacecraft configurations will be considered: a pair of spacecraft with X-band DOR tones (with a spanned bandwidth of $\Delta\nu_{DOR} = 40$ MHz) or a pair of spacecraft with Ka-band DOR tones ($\Delta\nu_{DOR} = 250$ MHz). (Additional DOR tones with smaller spanned bandwidths could be added to the spacecraft downlink spectrum, as required, to enable reliable ambiguity resolution.)

The mean elevation of the two spacecraft at each station will be assumed to be 20 deg, and the angular separation between the two spacecraft will be assumed to be solely in the elevation direction, providing a worst-case estimate of propagation media errors. Table 1 summarizes the observation description. This article considers spacecraft-spacecraft angular separations ranging from 0-20 deg on the sky plane.

At each station, two (or more) DOR tones from each spacecraft must be tracked to form the spacecraft DOR group delay observables. A multichannel closed-loop digital tracking receiver will simultaneously phase track all the tones from both spacecraft. While this capability does not currently exist in NASA's operational Deep Space Network (DSN), a demonstration is underway to use modified Global Positioning System (GPS) digital tracking receivers to simultaneously track carrier tones from Pioneer Venus Orbiter and Magellan at two stations [2]. Planned upgrades of the operational VLBI system will incorporate this capability. The closed-loop tracking provides a much lower data rate relative to current open-loop recording, and the multichannel capability provides higher SNR by eliminating time-multiplexing among DOR channels.

The actual spacecraft-spacecraft Δ DOR observable is obtained by combining the measured phases as follows: Let ν_{ij} represent the frequency of the i th DOR tone for spacecraft j , and let ϕ_{ijk} be the measured phase for that DOR tone at station k . The single-station one-way delay for spacecraft j at station k can then be obtained from a pair of DOR tones

$$\tau_{jk} = \frac{\phi_{2jk} - \phi_{1jk}}{\nu_{2j} - \nu_{1j}}$$

This one-way delay contains a bias due to uncertainty in the time of transmission of the signal from the spacecraft. By differencing this one-way observable between

two ground stations, this bias is eliminated. The unbiased, station-differenced delay observable for spacecraft j is then

$$\tau_j = \tau_{j2} - \tau_{j1}$$

Interpolating the two observations of spacecraft A (at $t = -T$ and $t = +T$, respectively) to the epoch of the observation of spacecraft B (at $t = 0$), and then differencing between spacecraft, yields the final spacecraft-spacecraft Δ DOR observable

$$\tau = \frac{1}{2}(\tau_A(-T) + \tau_A(T)) - \tau_B(0)$$

(The simple arithmetic mean of the two observations of spacecraft A is appropriate in the absence of significant angular accelerations for spacecraft A. A more general interpolation scheme could be used to account for any large accelerations.)

This observable represents a measure of the geometric delay τ_g , which is a function of the relative angular position of the two spacecraft

$$\tau_g = \frac{1}{c} \vec{B} \cdot (\hat{s}_A - \hat{s}_B)$$

where \hat{s}_A and \hat{s}_B are the unit vectors in the directions of the two spacecraft, \vec{B} is the baseline vector between ground antennas, and c is the speed of light.

In the next section, various error sources which corrupt this spacecraft-spacecraft Δ DOR observable will be examined. Each error source will be characterized in units that are most natural for the physical source of error, but ultimately one is interested in the angular error incurred on DSN intercontinental baselines. The following conversion factors will be used to relate various physical errors to an angular error on the sky plane:

$$\begin{aligned} 1\text{-cm path delay error} &= 33\text{-psec delay error} \\ &= 1.67\text{-nrad angular error} \end{aligned}$$

(This assumes a 6000-km projected length of the DSN intercontinental baseline on the sky plane.)

B. Error Components

1. **Spacecraft Signal-to-Noise Ratio.** The phase error σ_ϕ in the determination of the spacecraft tone phase is related to the SNR of the received DOR tone. The received DOR tone power can be expressed

$$P_{DOR} = P_{S/C} \zeta_{DOR} g_{S/C} \frac{\lambda^2}{(4\pi R)^2} g_{DSN}$$

where

$P_{S/C}$ = total transmitted spacecraft power

ζ_{DOR} = fraction of spacecraft power in the DOR tone (depends on modulation index and telemetry status)

$g_{S/C}$ = spacecraft antenna gain

λ = RF wavelength

R = Earth-spacecraft range

g_{DSN} = DSN ground antenna gain

In considering DOR tone SNR, the Mars Observer spacecraft will be used as a strawman configuration.² At maximum Earth-Mars range, with telemetry on (with an 80-deg modulation index), a 34-m high-efficiency DSN antenna provides a received DOR tone power of $P_{DOR} = -159.0$ dBm. The noise power per unit bandwidth is given by kT_{sys} , where k is Boltzmann's constant and T_{sys} is the system temperature of the receiving system, which represents the sum of the noise temperature of the first-stage amplifier, the brightness temperature of the atmosphere in the direction of the spacecraft, the 2.7-K cosmic background radiation, and any ground pickup from antenna spillover. Assuming a total noise system temperature of 25 K at X-band yields an X-band noise power per unit bandwidth of -184.6 dBm/Hz. The ratio of P_{DOR} to kT_{sys} , which describes the achievable link SNR in a one-second integration, is thus 25.6 dB-Hz.

The thermal phase error on the measured DOR tone phase is then given roughly by

$$\sigma_\phi = \sqrt{\frac{kT_{sys}}{P_{DOR} 2\tau_{int}}} \text{ rad}$$

where τ_{int} is the integration time of the observation. For the 120-sec integrations for each spacecraft assumed here,

² Ibid.

one arrives at a phase error of 5.4×10^{-4} cycles at X-band. (For the purposes of treating the statistical error due to SNR, one can treat the two 60-sec observations of spacecraft A as a single 120-sec scan at the same epoch as the observation of spacecraft B.) The final thermal delay error is thus given by

$$\sigma_\tau = \sqrt{2 \times 2 \times 2} \frac{\sigma_\phi}{\Delta\nu_{DOR}}$$

where the three factors of $\sqrt{2}$ reflect the pairwise differencing between DOR tones, stations, and spacecraft, resulting in an X-band delay error of 38.3 psec.

If one assumes a similar ratio of P_{DOR} to T_{sys} at Ka-band (which provides a reasonable figure of merit in designing the Ka-band DOR transponder), then one obtains an equivalent phase error for Ka-band VLBI. The resulting delay error would then be 6.1 psec due to the larger Ka-band spanned bandwidth.

2. **Ground System Instrumental Dispersion.** Uncalibrated phase dispersion in the ground receiving instrumentation induces errors in the measured tone phases that will corrupt the final spacecraft-spacecraft Δ DOR observable. As other error sources are reduced due to high SNR and common-mode cancellation of media effects, these dispersive errors may well represent a limiting error source for spacecraft-spacecraft Δ DOR. With current VLBI instrumentation, preliminary studies indicate that dispersive errors are at the 1- to 2-deg level,³ although more data on this error source are sorely needed. Achieving this level of phase error requires the use of phase calibration tones and/or the careful selection of a baseband frequency configuration to cancel instrumental errors between DOR tone channels.

A next-generation VLBI system employing broadband digitization of the entire intermediate frequency bandwidth and digital baseband filtering could significantly reduce instrumental errors by eliminating the analog baseband components that currently generate much of the dispersive phase effects. Design goals for this system provide for a one-millicycle dispersive phase error. The authors take this as the assumed instrumental dispersive phase error for each tone phase measurement. The resulting error in the spacecraft-spacecraft Δ DOR delay observable is

³ C. D. Edwards and K. Zukor, "Video Converter Local Oscillator Stability for Block I and Block II VLBI," JPL Interoffice Memorandum 335.1-90-055 (internal document), Jet Propulsion Laboratory, Pasadena, California, October 30, 1990.

$$\sigma_{\tau} = \sqrt{2 \times 2 \times 2} \frac{\sigma_{\phi}^{inst}}{\Delta\nu_{DOR}}$$

where the three factors of $\sqrt{2}$ again account for the pairwise differencing between DOR tones, stations, and spacecraft. This points out an important advantage of the increased spanned bandwidth available at Ka-band: For a given level of phase dispersion, delay errors are reduced proportional to the DOR tone spanned bandwidth. Assuming that $\sigma_{\phi}^{inst} = 1$ mcyc yields a σ_{τ} of 70.7 psec at X-band and 11.3 psec at Ka-band. At the smallest spacecraft angular separations, this error source will be one of the dominant contributors to the spacecraft-spacecraft Δ DOR error budget.

3. Station Clock Stability. Here the term "clock stability" represents both the stability of the station clock reference and the stability of the station frequency and timing distribution systems. The group delay error due to clock instability is on the order of

$$\Delta\tau = \sqrt{2} \times \sigma_y(\tau = 150 \text{ sec}) \times 150 \text{ sec}$$

where 150 sec is the time between central epochs of the scans for spacecraft A and B, and $\sigma_y(\tau = 150 \text{ sec})$ is the Allan standard deviation, or fractional frequency stability, evaluated at this time separation. Assuming a station stability of $\sigma_y(\tau = 150 \text{ sec}) = 10^{-14}$, this yields a spacecraft-spacecraft Δ DOR delay error of 2.1 psec. For a flicker-frequency noise spectrum, this error will grow linearly with the temporal scan separation.

4. Troposphere. The troposphere error can be separated into a static component and a fluctuating component. The static component represents the error made in the context of a static, isotropic refractivity distribution characterized by a single zenith troposphere delay. The delay at an arbitrary elevation angle θ is related to this zenith value by a mapping function f_{map} that is approximated here as $1/\sin\theta$. At a single station, an error σ_{τ}^{zen} in the zenith troposphere will lead to a delay error when differencing between spacecraft

$$\sigma_{\tau} = \sigma_{\tau}^{zen} \left| \frac{1}{\sin\theta_A} - \frac{1}{\sin\theta_B} \right|$$

where θ_i is the elevation angle of spacecraft i . For the DSN stations, σ_{τ}^{zen} is currently about 4 cm, based on seasonal weather models and surface meteorology. Water vapor radiometers and/or global GPS tracking data should

be able to provide reliable one-centimeter zenith troposphere estimates in the mid-1990s [3,4]; one centimeter will be used here as the representative zenith delay error. For two spacecraft with a mean elevation angle of 20 deg and angular separation $\Delta\theta$, assumed to be fully in the elevation direction, and accounting for uncorrelated one-centimeter zenith troposphere errors at each station, the resulting spacecraft-spacecraft Δ DOR delay error is

$$\sigma_{\tau} = \sqrt{2} \times 1 \text{ cm} \times \left| \frac{1}{\sin(20 \text{ deg} + \Delta\theta/2)} - \frac{1}{\sin(20 \text{ deg} - \Delta\theta/2)} \right|$$

In fact, the troposphere is neither static nor isotropic; spatial and temporal fluctuations, particularly in the distribution of atmospheric water vapor, lead to additional errors. Treuhaft and Lanyi [5] have developed a model of these fluctuations that is based on Kolmogorov turbulence. This model has been used to calculate the expected additional fluctuation error for the A-B-A scan sequence considered here, with the scans at a mean elevation of 20 deg and separated by 150 sec. The authors assume a tropospheric scale height of one kilometer, a wind speed of 8 m/sec, and a turbulence normalization constant of $2.4 \times 10^{-7} \text{ m}^{-1/3}$ [5]. For a zero-degree angular separation, the effect of temporal fluctuations over the 150-sec scan separation times yields a fluctuation error of about 10 psec; as the angular separation is increased, the additional effect of spatial fluctuations becomes important, with the total fluctuation error reaching about 39 psec for a 20-deg angular separation.

While not assumed in this analysis, it should be mentioned that improved line-of-sight troposphere calibrations (using either improved WVRs or lidar calibration techniques) could ultimately reduce the total wet troposphere error to well below one centimeter, independent of angular separation.

5. Ionosphere. Dual-frequency downlinks on both spacecraft would enable charged particle-induced errors to be virtually eliminated from the spacecraft-spacecraft Δ DOR delay observable. For the analysis here, however, the authors assume only a single-band downlink and calculate the size of the charged particle error that is incurred. The total ionospheric delay along a given line of sight can be expressed

$$\tau_{[\text{psec}]}^{ion} = 1340 \times \frac{TEC_{[10^{16} \text{ el/m}^2]}}{\nu_{[\text{GHz}]^2}}$$

where TEC is the line-of-sight integrated total electron content, and all units are indicated [6]. As with the troposphere, one can also separate the ionospheric error into a "static" and a "fluctuating" component. The mapping function f_{map} used to express the elevation dependence of the static component differs slightly from the tropospheric mapping function, due to the height of the ionospheric shell above the Earth, and takes the form

$$f_{map} = 1/\sin\left(\cos^{-1}\left[\frac{\cos\theta}{1+h/R}\right]\right)$$

where h is the height of the ionospheric shell above the Earth (~ 350 km) and R is the Earth's radius (~ 6371 km). The main impact is that the ionospheric mapping function increases more slowly at low elevations, saturating at a value of about 3.1 at the horizon. (This is a highly simplified picture of the ionosphere; in practice, the mapping function used is more complicated and accounts for the position of the Sun relative to the desired line of sight and the line of sight at which the ionospheric calibration was performed, in order to account for the diurnal variation in TEC . Nevertheless, the simple picture used here is adequate to estimate a typical error gradient on the sky due to ionospheric calibration error.) Using an analysis similar to that used for the troposphere shows that the static ionospheric error for the spacecraft-spacecraft ΔDOR delay observable is

$$\sigma_{\tau[\text{psec}]} = \sqrt{2} \times 1340 \times \frac{5}{\nu_{[\text{GHz}]}} \times |f_{map}(20 \text{ deg} + \Delta\theta/2) - f_{map}(20 \text{ deg} - \Delta\theta/2)|$$

where the authors have assumed an uncertainty in the zenith TEC of $\sigma_{TEC} = 5 \times 10^{16}$ el/m², and where the authors again take the worst-case geometry for which the spacecraft angular separation is entirely in the elevation direction. This corresponds to a 10-percent calibration uncertainty for a typical daytime maximum of $TEC = 50 \times 10^{16}$ el/m², which is consistent with ionospheric calibration accuracies using Faraday rotation or GPS satellite data. For $\Delta\theta = 5$ deg, this represents a delay error of 36 psec at X-band, or 2 psec at Ka-band.

The fluctuating component for the ionosphere is expected to be important, due to the variety of phenomena driving the ionospheric charged particle distribution (e.g., the day-night asymmetry, traveling ionospheric disturbances, and latitude variations) and the resulting limited accuracy of the simple static ionosphere model. Theoretical understanding of the processes driving ionospheric

fluctuations is much less developed than for tropospheric fluctuations; as a result, empirical data will be used to guide the quantitative estimate of this error source. Based on a recent study that derives temporal fluctuation statistics from dual-frequency GPS carrier phase data,⁴ an additional error of 0.5 TEC units (1 TEC unit = 10^{16} el/m²) is specified to account for temporal ionospheric fluctuations at each site on the time scale of the differential observations. This level of fluctuation corresponds to a delay error of 13 psec at X-band and 0.9 psec at Ka-band.

6. Solar Plasma. Charged particles in the solar plasma also induce a delay error for spacecraft-spacecraft ΔDOR . The solar plasma delay error is proportional to the double-differenced line-of-sight integrated total electron content in the solar wind along the four relevant spacecraft-spacecraft ΔDOR lines of sight. The statistical model of Kahn and Border [7], based on solar plasma electron density spectra compiled by Woo and Armstrong [8], derives a spatial structure function for solar plasma-induced phase fluctuations, which can be used to calculate the error in the station-differenced delay to a single spacecraft

$$\sigma_{\tau} = \frac{134 \text{ psec}}{\nu_{[\text{GHz}]}^2 [\sin SEP]^{1.225}}$$

where SEP is the Sun-Earth-probe angle. For spacecraft-spacecraft angular separations greater than about one degree, the solar plasma error for a second spacecraft will be essentially uncorrelated. Assuming a projected DSN baseline of 6000 km then yields a total spacecraft-spacecraft ΔDOR error of

$$\sigma_{\theta} = \frac{9.50 \text{ nrad}}{\nu_{[\text{GHz}]}^2 [\sin SEP]^{1.225}}$$

A SEP angle of 20 deg for both spacecraft is assumed, which yields a total angular error of $\sigma_{\theta} = 0.50$ nrad at X-band and 0.03 nrad at Ka-band. Below a one-degree angular separation, the solar plasma error will be further reduced due to additional cancellation between raypaths for the two spacecraft.

7. Baseline Errors. Uncertainty in the baseline vector is due to a combination of a priori station location errors and errors in the knowledge of Earth orientation

⁴ A. J. Mannucci, "Temporal Statistics of the Ionosphere," JPL Interoffice Memorandum 335.1-90-056 (internal document), Jet Propulsion Laboratory, Pasadena, California, October 25, 1990.

(UT1-UTC and polar motion). Any uncertainty $\delta\vec{B}$ in the baseline vector leads to a delay error

$$\delta\tau = \frac{1}{c} \delta\vec{B} \cdot (\hat{s}_1 - \hat{s}_2)$$

where \hat{s}_1 and \hat{s}_2 are the source directions to the two spacecraft. (In other words, the baseline path delay error is attenuated by the angular spacecraft separation, in radians, projected along the baseline direction.) It is assumed here that station coordinates in the terrestrial frame are known to 3 cm per component [9]. In addition, weekly VLBI observations combined with daily GPS observations have been shown to be able to deliver real-time Earth orientation estimates with 10-nrad accuracy [10]. Based on these two error components, a delay error of $\sigma_\tau = 4.3$ psec $\times \Delta\theta_{[\text{deg}]}$ due to baseline uncertainty is specified.

8. Frame Tie. One final error contribution is related to the offset in the planetary and radio reference frames, and is referred to as the frame-tie uncertainty. In conventional spacecraft-quasar Δ VLBI, the spacecraft position is measured in the radio frame relative to a nearby quasar; the frame-tie offset contributes directly as an angular bias for determining the spacecraft position relative to a planetary target. The spacecraft-spacecraft Δ VLBI technique described in this article reduces the effect of the frame-tie error by directly measuring the approach spacecraft relative to a spacecraft at the target planet. Nonetheless, the frame tie does induce a small residual error in converting the measured Δ VLBI delay into an angular separation. The error is due to the fact that the baseline orientation is modeled in the radio reference frame, based on periodic VLBI and GPS measurements of Earth orientation, while the reference spacecraft's position is tied to the planetary ephemeris. The resulting angular error in the approach spacecraft's angular position relative to the target planet is proportional to the product of the frame-tie uncertainty and the angular separation between spacecraft, expressed in radians. The frame-tie uncertainty is currently about 50 nrad for the inner planets, but that value should be reduced to about 25 nrad based on several ongoing observational programs, including millisecond pulsar timing and VLBI observations [11], observations of planetary occultations of quasars [12], and joint solutions of VLBI and lunar laser ranging data sets [13]. Thus, the authors include an error in the determination of the approach spacecraft's target-relative position in terms of the spacecraft-spacecraft angular separation $\delta\theta$

$$\sigma_\theta = 25 \text{ nrad} \times \frac{\pi}{180} \times \delta\theta_{[\text{deg}]} = 0.44 \text{ nrad} \times \delta\theta_{[\text{deg}]}$$

C. The Total Spacecraft-Spacecraft Δ DOR Error Budget

Table 2 summarizes the error-modeling assumptions made in this analysis, while Tables 3 and 4 present the error budget for the X-band and Ka-band spacecraft-spacecraft Δ DOR cases considered here. Figure 3 summarizes the angular error for each case as a function of the angular separation between spacecraft. For the X-band case, the dominant errors for large angular separations (>10 deg) are the propagation media errors, due to uncertainties in the zenith troposphere and ionosphere delays. The angular error grows roughly by 0.7 nrad per degree of angular separation in this range. For smaller angular separations, the dominant errors are the small-scale fluctuations in the ionosphere and the instrumental phase dispersion, followed by troposphere fluctuations and the statistical measurement error due to the received spacecraft SNR. As the angular separation approaches zero, the accuracy levels out at just over 4 nrad.

The Ka-band error budget shows further accuracy improvement due to two factors. First, the much larger spanned bandwidth reduces the statistical measurement error as well as the phase dispersion error by a factor of 250/40 relative to the X-band case. Second, the higher Ka-band frequency reduces the effects of the ionosphere and the solar plasma by a factor of $(32/8.4)^2$, or about 14.5. For Ka-band, the dominant errors are troposphere and platform errors at large angular separations, and instrumentation and troposphere fluctuations at small angular separations.

III. Discussion

The error budget presented in the last section was parameterized as a function of the angular separation of the approach and in-orbit spacecraft. How does this angular separation evolve during the final weeks of planetary approach? As a representative example, consider the spacecraft-spacecraft angular separation for a Hohmann (minimum-energy) Earth-Mars transfer orbit. For this orbit, the Mars-Earth-probe (MEP) angle is less than 27 deg for essentially the entire trajectory, less than 10 deg for the last four months of the trajectory, and, in fact, less than 2 deg for the last 100 days. At encounter, the rate of change of the MEP angle is only 0.044 deg/day.

Higher energy transfer orbits, for which aerocapture insertions might be a key component, and therefore which may require highest accuracy approach navigation, will typically have larger approach velocities, but should still

have a spacecraft-planet angular separation of less than 5 deg for at least the last few weeks of planetary approach. Hence, the highest accuracy spacecraft-spacecraft Δ DOR observables will be available during the final critical targeting maneuvers in the last few weeks prior to encounter.

It should also be mentioned that during the final hours of planetary approach, the in-orbit and approach spacecraft will become sufficiently close on the sky plane that they may be observed simultaneously within a single Earth-based antenna beamwidth. This enables the use of the same-beamwidth interferometry (SBI) technique [2,14], in which the simultaneous observation of both spacecraft leads to further significant error reductions, with accuracies of 10-100 prad possible if the RF phase observable can be resolved. The X-band and Ka-band beamwidths of a 34-m antenna are 60 and 16 mdeg, respectively; thus, for the Hohmann trajectory described above, X-band SBI observations will be possible for over a day before encounter, and Ka-band for about 8 hours prior to encounter. Konopliv and Wood [15] have already shown how the SBI observable can provide accurate Mars approach navigation during the final hours of approach, helping to enable aerocapture. The key message of the results presented here, however, is that even before SBI observations are possible, nonsimultaneous spacecraft-spacecraft Δ DOR observations can provide significant improvements in target-relative approach navigation for weeks or even months before encounter.

One important error source that applies to spacecraft-spacecraft Δ DOR was not included in the error budget presented here: namely, the uncertainty in the planet-relative position of the in-orbit reference spacecraft. This error depends very much on the type of orbit the reference spacecraft is in, the amount and quality of tracking data collected for the in-orbit spacecraft, and assumptions about limiting errors, such as uncertainties in the planetary gravity field. Preliminary navigation analysis for the Mars Observer mission, for example, indicates that one-kilometer orbit errors are expected immediately after orbit insertion. However, after several weeks of intensive Doppler tracking, the resulting improvement in the Mars gravity field should allow a reduction of orbit errors to about 200 meters.⁵ A 200-m spacecraft position error corresponds to an angular error of 0.5-2.5 nrad, depending on the Earth-Mars range. For X-band observations, this error will not be dominant, but it will be an important error source for the higher accuracy Ka-band observations.

⁵ P. Esposito, S. Demcak, D. Roth, G. Bollman, and A. Halsell, "Mars Observer Project Navigation Plan," JPL D-3820 (internal document), Jet Propulsion Laboratory, Pasadena, California, June 15, 1990.

Of course, if the reference spacecraft is, in fact, a beacon on the planetary surface, its position will be known to a much higher accuracy: Conventional range and Doppler data should be able to provide Mars-centered beacon position determination with 10-m accuracy in the spin radius and 100-m accuracy along the spin axis. In addition, SBI between the surface beacon and an orbiter could provide few-meter Mars-relative beacon position accuracy in all three components [16].

Some error sources which are important for conventional quasar-relative Δ DOR are eliminated or greatly reduced in spacecraft-only observations. In the previous section the authors discussed how the frame-tie error is greatly reduced in the spacecraft-spacecraft technique, relative to spacecraft-quasar Δ VLBI. Other important error sources for quasar-relative Δ DOR include the statistical uncertainty in the measurement of the quasar delay, any a priori uncertainty in the quasar position, and the effect of source structure on the quasar position. Because of the limited 250-KHz recorded bandwidth of the NCB VLBI system, the quasar delay measurement error is one of the limiting error sources for conventional Δ DOR. An additional impact of the low NCB sensitivity is that only bright quasars can be reliably observed: A minimum correlated flux density of 0.4 Jy is typically required for reliable detection with a pair of DSN antennas, one 70-m and one 34-m. Due to the limited number of useful sources, it is often necessary to use a quasar more than 10 deg from the spacecraft, with the result that Earth orientation and propagation media errors are increased.

Finally, a priori source positions are uncertain at the level of about 5 nrad, based on the current DSN quasar data set. Source position accuracies may improve further, toward one-nanoradian accuracy with the increasing amount of Mark III observations in the source catalog data set. However, it is suspected that source structure can cause few-nanoradian errors in apparent source position, varying with time as the quasar jet structure evolves, and also changing with observation geometry as the fringe orientation changes and different source features are resolved; this error source may pose a difficult obstacle to achieving nanoradian-level quasar positions. Over short periods of weeks or months, during which source structure is expected to remain nearly constant, it is possible to eliminate the source structure error in a relative sense among a series of Δ DOR observations by always observing the same quasar(s) at exactly the same hour angle(s) [1]. However, this poses scheduling constraints and still does not eliminate any overall source position error common to all observations.

IV. Opportunities to Demonstrate Spacecraft-Spacecraft Δ DOR

To validate the error budget presented here and gain experience in acquiring and processing this data type, it would be valuable to find opportunities to demonstrate the spacecraft-spacecraft Δ DOR observation technique. To demonstrate the technique, one requires two angularly close spacecraft, each with downlinks at the same frequency band, including VLBI DOR tones. Two noteworthy opportunities are mentioned here. First, in January of 1994, Venus and Mars pass near each other on the sky plane. At this time, Magellan will be in orbit about Venus, and Mars Observer will have recently arrived at Mars. Mars Observer has a 38.25-MHz DOR tone bandwidth, and Magellan has a 30.72-MHz bandwidth consisting of the ± 16 th harmonics of the 960-KHz telemetry subcarrier. The two spacecraft will pass within about 0.3 deg of each other, with closest approach occurring on January 6, 1994. By acquiring a series of spacecraft-spacecraft Δ DOR observations over a several-week period, it should be possible to verify the accuracy of the spacecraft-spacecraft Δ DOR data type as a function of spacecraft angular separation. In addition, a subkilometer determination of the sky-plane components of the offset between Venus and Mars at this epoch would result, providing a valuable constraint on the relative orientation of the orbit planes of these two planets.

The second, and potentially more interesting, opportunity involves using the Mars Observer spacecraft as a spacecraft-spacecraft Δ DOR reference for planetary approach of the Russian Mars '94 spacecraft. Mars '94 tentatively plans a September 1995 Mars orbit insertion; this is near the end of the prime mission of Mars Observer, which will have arrived at Mars in August of 1993. Mars Observer incorporates a 38.25-MHz X-band DOR tone bandwidth; if the Russians incorporate a similar capability on their spacecraft, it would be possible to collect spacecraft-spacecraft Δ DOR data during Mars '94 approach. In Part II of this article, a covariance analysis will be presented to examine the navigation benefits of such an observation program for the Russian spacecraft. And of course, after encounter, one would also be interested in collecting SBI data, which have already been shown to have significant navigational benefits to both missions [2,14]. This scenario provides a unique opportunity to demonstrate multiple spacecraft tracking at Mars, where future,

more ambitious aerocapture missions encompassed within the SEI will benefit from new, high-accuracy tracking techniques.

V. Summary

Spacecraft-spacecraft Δ DOR observations between an in-orbit spacecraft and another spacecraft approaching that planet can provide target-relative angular navigation with accuracies of about 4 nrad during the last weeks of planetary approach for spacecraft equipped with X-band transponders incorporating roughly 40-MHz DOR tone spacings. Accuracies approaching one nanoradian can be obtained by going to Ka-band downlinks with DOR tone spacings of several hundred megahertz. These accuracies correspond to an observation duration of only 6 minutes.

The spacecraft-spacecraft Δ DOR observable has the advantage of tying the approach spacecraft directly to the planetary target. In addition, because only spacecraft signals are used, no wideband quasar recording is required. As a result, data transfer and data processing are simplified, which enables these observables to be available in near-real-time. To enable efficient data collection, a key part of the DSN's planned VLBI system upgrade should be the implementation of ground tracking receivers that can simultaneously track multiple tones from each spacecraft. Reducing instrumental phase dispersion errors to the millicycle level will be an important design goal for this new system and should be achievable by using digital data acquisition techniques. Other key improvements in ground capabilities assumed in this analysis are a 1-cm zenith troposphere calibration capability and a 3-cm station location knowledge. In the 1995 time frame, GPS and/or WVRs should be capable of providing troposphere calibrations at this level, while VLBI, GPS, and LLR data should be able to provide the required level of station location accuracy.

Opportunities to demonstrate the spacecraft-spacecraft Δ DOR technique will arise in the next several years, first with fortuitous sky-plane flybys of unrelated deep space missions, such as Magellan and Mars Observer, and then in 1995 by the possibility of using differential observations of Mars Observer and Mars '94 to improve the planetary approach targeting for the Mars '94 mission.

Acknowledgments

This analysis benefits from earlier analyses of conventional spacecraft-quasar Δ DOR by Brooks Thomas, Sien Wu, and Bob Treuhaft. The authors thank Roger Linfield, Bill Folkner, Sam Thurman, and Lincoln Wood for their helpful comments on a draft of this article.

References

- [1] R. N. Treuhaft and S. T. Lowe, "Nanoradian VLBI Tracking for Deep Space Navigation," paper AIAA 90-2939, *Proceedings of the AIAA/AAS Astrodynamics Conference*, pp. 587-589, Part 2, Portland, Oregon, August 20-22, 1990.
- [2] W. Folkner and J. Border, "Orbiter-Orbiter and Orbiter-Lander Tracking Using Common-Beam Interferometry," paper 90-2906, *Proceedings of the AIAA/AAS Astrodynamics Conference*, Part 1, pp. 355-363, Portland, Oregon, August 20-22, 1990.
- [3] D. M. Tralli and S. M. Lichten, "Stochastic Estimation of Tropospheric Path Delays in Global Positioning System Geodetic Measurements," *Bulletin Geodesique*, vol. 64, pp. 127-159, 1990.
- [4] T. H. Dixon and W. S. Kornreich, "Some Tests of Wet Tropospheric Calibration for the CASA Uno Global Positioning System Experiment," *Geophys. Res. Letters*, vol. 17, pp. 203-206, March 1990.
- [5] R. N. Treuhaft and G. E. Lanyi, "The Effect of the Dynamic Wet Troposphere on Radio Interferometric Measurements," *Radio Sci.*, vol. 22, pp. 251-265, 1987.
- [6] P. S. Callahan, "Ionospheric Variations Affecting Altimeter Measurements: A Brief Synopsis," *Marine Geodesy*, vol. 8, pp. 249-263, 1984.
- [7] R. D. Kahn and J. S. Border, "Precise Interferometric Tracking of Spacecraft at Low Sun-Earth-Probe Angles," paper AIAA-88-0572, presented at Aerospace Sciences Meeting, Reno, Nevada, January 11-14, 1988.
- [8] R. Woo and J. W. Armstrong, "Spacecraft Radio Scattering Observations of the Power Spectrum of Electron Density Fluctuations in the Solar Wind," *J. Geophys. Res.*, vol. 84, p. 7288, 1979.
- [9] R. P. Malla and S. C. Wu, "GPS Inferred Geocentric Reference Frame for Satellite Positioning and Navigation," *Bulletin Geodesique*, vol. 63, pp. 263-279, 1989.
- [10] A. P. Freedman, "Determination of Earth Orientation Using the Global Positioning System," *TDA Progress Report 42-99*, vol. July-September 1989, pp. 1-11, November 15, 1989.
- [11] D. Jones, R. Dewey, C. Gwinn, and M. Davis, "Mark III VLBI Astrometry of Pulsars," *IAU Colloquium 131, Radio Interferometry: Theory, Techniques, and Applications*, ed. T. Cornwell, Socorro, New Mexico, October 1990.
- [12] R. Linfield, "Using Planetary Occultations of Radio Sources for Frame Tie Measurements; Part 1: Motivation and Search for Events," *TDA Progress Report 42-103*, vol. July-September 1990, pp. 1-13, November 15, 1990.

- [13] M. Finger and W. Folkner, "A Determination of the Radio-Planetary Frame-Tie and the DSN Tracking Station Locations," paper 90-2905, *Proceedings of the AIAA/AAS Astrodynamics Conference, Part 1*, Portland, Oregon, pp. 335-353, August 20-22, 1990.
- [14] J. Border and W. M. Folkner, "Differential Spacecraft Tracking by Interferometry," paper CNES-89-145, *Proceedings of the CNES International Symposium on Space Dynamics*, Toulouse, France, pp. 6-10, November 1989.
- [15] A. Konopliv and L. Wood, "High Accuracy Mars Approach Navigation with Radiometric and Optical Data," paper AIAA 90-2907, *Proceedings of the AIAA/AAS Astrodynamics Conference, Part 1*, Portland, Oregon, pp. 364-376, August 20-22, 1990.
- [16] R. D. Kahn, W. M. Folkner, C. D. Edwards, and A. Vijayaraghavan, "Position Determination of Spacecraft at Mars Using Earth-Based Differential Tracking," paper AAS 91-502, presented at AAS/AIAA Astrodynamics Specialist Conference, Durango, Colorado, August 19-22, 1991.

Table 1. Observation description.

Observation sequence	Time, sec	
Spacecraft A	60	
Slew time	60	
Spacecraft B	120	
Slew time	60	
Spacecraft A	60	
Observation geometry		
Mean elevation angle	20 deg	
Angular separation	0-20 deg, in elevation direction at both stations	
Projected baseline length	6,000 km	
Spacecraft signal spectrum		
	Case 1, X-band	Case 2, Ka-band
Carrier frequency	8.4 GHz	32.0 GHz
DOR tone spacing	±20 MHz	±125 MHz
Received DOR tone SNR	25.6 dB-Hz	25.6 dB-Hz

Table 2. Error-modeling assumptions.

Spacecraft SNR	
DOR bandwidth	
Case 1, X-band	40 MHz
Case 2, Ka-band	250 MHz
P_{tone}/N_0	25.57 dB-Hz
Instrumentation	
Single-channel dispersive phase error	0.001 cyc
Clock stability	
Time between spacecraft scans	150 sec
Allan variance	1×10^{-14}
Static troposphere	
Zenith troposphere uncertainty	1 cm
Mean elevation angle	20 deg
Fluctuating troposphere	
Treuhaft-Lanyi model	(as per [8])
Static ionosphere	
Zenith ionosphere uncertainty	5 TEC units
Frequency	
Case 1, X-band	8.4 GHz
Case 2, Ka-band	32 GHz
Mean elevation angle	20 deg
Fluctuating ionosphere	
RMS TEC fluctuation	0.5 TEC units
Baseline	
Station location uncertainty	3 cm
Earth orientation uncertainty	10 nrad
Radio-planetary frame tie	
Frame-tie error	25 nrad
Solar plasma	
Sun-Earth-probe angle	20 deg

Table 3. Error budget for X-band spacecraft-spacecraft Δ DOR.

Angular separation, deg	Spacecraft SNR, nrad	Instrumentation, nrad	Clock stability, nrad	Troposphere, nrad	Ionosphere, nrad	Base-line, nrad	Solar plasma, nrad	Frame tie, nrad	RSS, nrad
0.10	1.91	3.54	0.11	0.67	0.67	0.02	<0.50	0.22	4.16
1.00	1.91	3.54	0.11	0.76	0.76	0.21	0.50	0.44	4.22
2.00	1.91	3.54	0.11	0.97	0.99	0.43	0.50	0.87	4.39
4.00	1.91	3.54	0.11	1.55	1.59	0.86	0.50	1.75	5.01
8.00	1.91	3.54	0.11	2.95	2.95	1.71	0.50	3.49	7.00
10.00	1.91	3.54	0.11	3.74	3.64	2.14	0.50	4.36	8.20
20.00	1.91	3.54	0.11	9.29	7.00	4.28	0.50	8.73	15.69

Table 4. Error budget for Ka-band spacecraft-spacecraft Δ DOR.

Angular separation, deg	Spacecraft SNR, nrad	Instrumentation, nrad	Clock stability, nrad	Troposphere, nrad	Ionosphere, nrad	Base-line, nrad	Solar plasma, nrad	Frame tie, nrad	RSS, nrad
0.10	0.31	0.57	0.11	0.67	0.05	0.02	<0.03	0.22	0.94
1.00	0.31	0.57	0.11	0.76	0.05	0.21	0.03	0.44	1.11
2.00	0.31	0.57	0.11	0.97	0.07	0.43	0.03	0.87	1.52
4.00	0.31	0.57	0.11	1.55	0.11	0.86	0.03	1.75	2.57
8.00	0.31	0.57	0.11	2.95	0.20	1.71	0.03	3.49	4.93
10.00	0.31	0.57	0.11	3.74	0.25	2.14	0.03	4.36	6.17
20.00	0.31	0.57	0.11	9.29	0.48	4.28	0.03	8.73	13.47

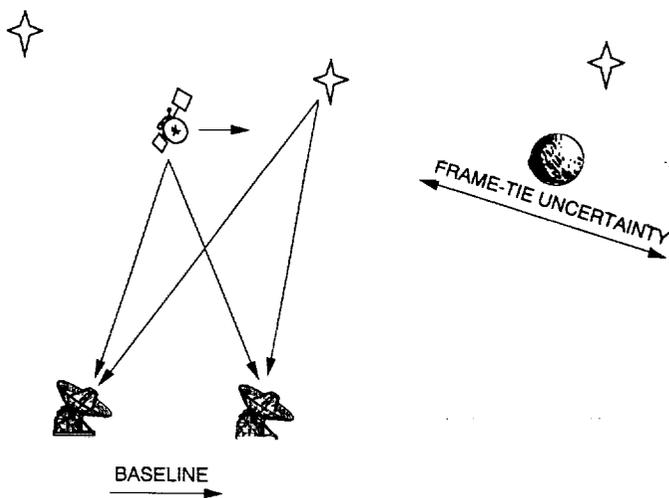


Fig. 1. Conventional Δ DOR provides a determination of the angular position of a spacecraft relative to the reference frame of distant quasars. Uncertainty in the position of the target planet in this reference frame represents an important navigation error source.

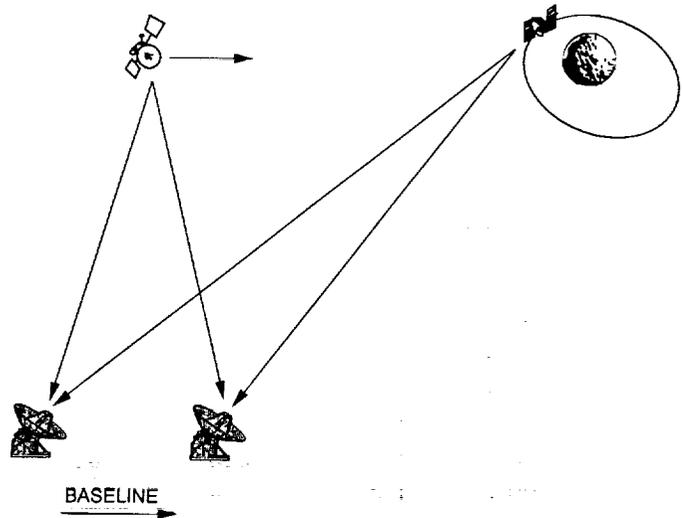


Fig. 2. Spacecraft-spacecraft Δ DOR observations between an approach spacecraft and a planetary orbiter provide direct planet-relative approach navigation.

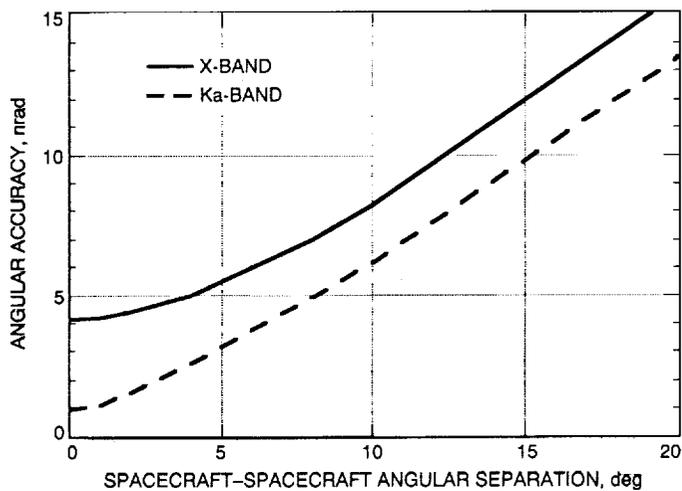


Fig. 3. Angular accuracy versus spacecraft angular separation for X-band and Ka-band spacecraft-spacecraft Δ DOR.

N 93 - 19419

128439

P-15

Use of the VLBI Delay Observable for Orbit Determination of Earth-Orbiting VLBI Satellites

J. S. Ulvestad
Navigation Systems Section

Very long-baseline interferometry (VLBI) observations using a radio telescope in Earth orbit were performed first in the 1980s. Two spacecraft dedicated to VLBI are scheduled for launch in 1995; the primary scientific goals of these missions will be astrophysical in nature. This article addresses the use of space VLBI delay data for the additional purpose of improving the orbit determination of the Earth-orbiting spacecraft. In an idealized case of quasi-simultaneous observations of three radio sources in orthogonal directions, analytical expressions are found for the instantaneous spacecraft position and its error. The typical position error is at least as large as the distance corresponding to the delay measurement accuracy but can be much greater for some geometries.

A number of practical considerations, such as system noise and imperfect calibrations, set bounds on the orbit-determination accuracy realistically achievable using space VLBI delay data. These effects limit the spacecraft position accuracy to at least 35 cm (and probably 3 m or more) for the first generation of dedicated space VLBI experiments. Even a 35-cm orbital accuracy would fail to provide global VLBI astrometry as accurate as ground-only VLBI. Recommended changes in future space VLBI missions are unlikely to make space VLBI competitive with ground-only VLBI in global astrometric measurements.

I. Introduction

Very long-baseline interferometry (VLBI) is a radio astronomy technique that achieves high angular resolution by means of the simultaneous recording of signals from artificial or natural radio sources at widely separated radio telescopes, and then cross-correlating those signals at a central processing facility [1]. This technique has been used with ground radio telescopes for about 25 years. Currently its uses include high-resolution imaging of radio sources, radio-source position measurements, and moni-

toring of Earth-rotation parameters and continental drift (e.g., [2]). The first VLBI experiments involving a space radio telescope along with the ground radio telescopes were performed between 1986 and 1988 [3-6]. In those experiments, a Tracking and Data Relay Satellite (TDRS) was used as the space radio telescope, while large telescopes in Australia and Japan were used on the ground.

Two dedicated space VLBI satellites, with radio telescopes 7-10 m in diameter, are scheduled for launch in

1995. The VSOP (VLBI Space Observatory Programme) satellite is being developed in Japan [7], while Radioastron [8] will be a product of the Russian space agency. Radioastron will operate in a highly elliptical orbit with a perigee height of about 3000 km and an apogee height of 80,000 km, while the VSOP orbit will have a perigee height of 1000 km and an apogee height of 20,000 km. Other missions that would be launched after the year 2000 have been studied, particularly the International VLBI Satellite [9]. The primary goals of all these missions are astrophysical, with models and images of radio-source morphologies being the most important scientific output. In this article, the additional use of space VLBI data for improved orbit determination of space VLBI satellites and for radio-source position measurements ("astrometry"), first suggested in [10], are explored. However, the performance in these areas is not the justification for the space missions, and the limitations that are discussed in this article in no way imply limitations in the ability of the spacecraft to achieve their primary goals.

In theory, the longer baselines available to space radio telescopes could provide enhanced astrometry of radio sources, but only if the knowledge of the baselines has accuracy comparable to that available for ground-ground baselines. This requires position accuracies of a few centimeters for the space telescope, 2-3 orders of magnitude better than the accuracy of tens of meters achievable for high Earth orbiters using conventional Doppler tracking [11]. One possible way to obtain such high accuracy is by the use of Global Positioning System (GPS) receivers on board the spacecraft [12], although their use would be restricted somewhat for orbits outside the GPS constellation (altitudes of 20,000 km). It also has been suggested that the VLBI data acquired on space-ground baselines could be used to achieve the desired accuracy [13,14]. However, studies of this subject often have ignored potentially important error sources or have not obtained the desired positional accuracy. Even if centimeter-level baseline accuracy could not be obtained using space VLBI data, it still is of interest to determine the improvement in orbit determination that might be achieved through the use of such data.

This article represents one aspect of an effort to analyze the potential of space VLBI data for improved orbit determination. It has two main goals. The first is an analysis of the information content and the smallest possible orbit-determination errors for delay measurements in a highly idealized scenario of space VLBI observations. This framework is intended to provide a limiting case that shows the best performance that might be achieved for a hypothetical VLBI satellite. The second goal is to investigate a

number of practical limitations to the accuracy achievable in a realistic observing scenario. The effect of these limitations on the ability to do astrometry using space VLBI is addressed, with emphasis on the application to the first generation of dedicated space VLBI satellites that will operate in the 1990s.

II. Assumptions and Definitions

It is assumed that VLBI observations are performed using one space-based telescope and one or more telescopes on the surface of the Earth. The space- and ground-based telescopes simultaneously observe the same set of radio sources. The VLBI data at the ground telescopes are recorded in the standard way on high-capacity videotapes. The radio-source data received by the space antenna are digitized and transmitted to a ground tracking station for recording. The clock at the tracking station that records the time of data reception at the spacecraft is initialized by means of a tone sent up from the tracking station and transponded by the spacecraft. Thereafter, the ground clock is driven by the VLBI bitstream, with the evolution of the clock error monitored and corrected by means of a two-way phase link between the tracking station and the spacecraft. The videotapes from the ground and space telescopes then are correlated at a central processing facility, where amplitude, relative phase, delay, and delay-rate observables are extracted for each observation. A spacecraft orbit reconstructed from two-way Doppler data (e.g., [11,15]) is assumed to be the a priori model used at the correlation facility. In this article, only the VLBI delay observable is considered. It is assumed that a single ground tracking station is used during the set of VLBI observations considered.

In order to visualize the delay observable and its relationship to spacecraft position and clock uncertainties, it is appropriate to define a Cartesian coordinate system whose origin is at the location of the ground tracking station in contact with the spacecraft. Instantaneous unit vectors of that coordinate system are \hat{i} and \hat{j} at right angles to each other in the Earth's equatorial plane, with \hat{k} directed toward the North Pole. This coordinate system is a mixture of the classical topocentric and equatorial coordinate systems and is selected because it simplifies the analysis done below in Section III. A ground radio telescope has position $\vec{r}_t = (x_t, y_t, z_t)$ relative to the tracking station, while the spacecraft position is $\vec{r} = (x, y, z)$, with the direction from the tracking station to the spacecraft indicated by the unit vector \hat{r} . The range between the tracking station and the spacecraft is r . Figure 1 shows the geometry for a space VLBI observation of a single radio source. The measured VLBI delay τ is given approximately by

$$\tau = \tau_g + \tau_c + \tau_p + \tau_i + \tau_s \quad (1)$$

In Eq. (1), τ_g is the geometric delay for the space-ground baseline, τ_c is the delay contribution caused by the experiment clocks, τ_p is the delay caused by propagation through the Earth's troposphere and ionosphere, τ_i is the instrumental delay, and τ_s is the delay caused by radio-source structure. Note that the propagation delays must be included because they affect the data received by the ground radio telescope, hence contributing to the delay measured on the space-ground baseline.

III. Error Analysis and Information Content in an Idealized Case

Consider a VLBI observation of a hypothetical radio source located on the x -axis. In the ideal case, assume that the radio source is a point source ($\tau_s = 0$), the propagation delays caused by the ionosphere and troposphere are known perfectly (τ_p equals the modeled propagation delay, $\tau_{p,m}$), and there is no instrumental delay ($\tau_i = 0$). Further, assume that all clocks on the ground are perfect and that there are no errors in the determinations of universal time, Earth rotation, or polar motion. In that case, the only contribution to τ_c comes from the spacecraft clock. Suppose that the best model of the spacecraft position is \vec{r}_m , the best model of the ground telescope location is $\vec{r}_{i,m}$, the model geometric delay is $\tau_{g,m}$, and the best model of the spacecraft clock delay is $\tau_{c,m}$. There are a number of possible ways to find the best model of the clock delay; details are not considered here, and the reader is referred to [16] for more discussion. If the delay is "tracked" in the correlator by subtracting the model delay τ_m , the measured residual delay $\tau_{r,1}$ will be given by

$$\begin{aligned} c\tau_{r,1} &\equiv (\tau - \tau_m) \\ &= c(\tau_g - \tau_{g,m}) + c(\tau_c - \tau_{c,m}) \\ &= (\vec{r} - \vec{r}_m) \cdot \hat{i} - (\vec{r}_i - \vec{r}_{i,m}) \cdot \hat{i} + c(\tau_c - \tau_{c,m}) \quad (2) \end{aligned}$$

where the speed of light is given by c . If it is assumed further that the locations of the ground telescope and tracking station are known perfectly (i.e., $\vec{r}_i = \vec{r}_{i,m}$), Eq. (2) reduces to

$$c\tau_{r,1} = (\vec{r} - \vec{r}_m) \cdot \hat{i} + c(\tau_c - \tau_{c,m}) \quad (3)$$

The magnitude of the delay contributed by the spacecraft clock is the time it takes to transmit a tone between the spacecraft and the tracking station. The clock epoch is initialized and monitored as summarized in Section II. Because the clock's evolution is monitored precisely by means of the two-way phase link, only the spacecraft clock delay τ_c at the initialization epoch need be considered here. The transponded tone received at the tracking station at the initialization time arrives after a delay given by the light travel time, so the initial spacecraft data are tagged with a time later than the time at which the VLBI data actually were received by the space radio telescope. This delay must be modeled in the correlation of the VLBI data. Since the clock delay at the initialization time t_0 is just r_0/c , where $r_0 \equiv r(t_0)$, Eq. (3) reduces to

$$c\tau_{r,1} = (x - x_m) + (r_0 - r_{0,m}) \quad (4)$$

The simple analysis above shows that the measured delay residual depends on the errors in the spacecraft location in two different (usually not orthogonal) directions. One is the line-of-sight direction to the infinitely distant radio source (\hat{i}), while the other is the line of sight from the tracking station to the spacecraft (\hat{r}_0). A single measurement of $\tau_{r,1}$ cannot distinguish between the two. Simultaneous observations from other ground radio telescopes cannot contribute new information. The location of a ground radio telescope contributes only to the term $(\vec{r}_i - \vec{r}_{i,m})$ in Eq. (2), but (by assumption) this term is negligible. One can use as many ground telescopes as are available on Earth, and the result still reduces to Eq. (4), one equation in two unknowns, x and r_0 .¹

In an ideal world, one could imagine making a VLBI observation instantaneously, then slewing to another source at infinite speed and making another instantaneous observation. If observations are made of radio sources in the \hat{j} and \hat{k} directions in this hypothetical world, the following equations are added to the system:

$$c\tau_{r,2} = (y - y_m) + (r_0 - r_{0,m}) \quad (5)$$

and

$$c\tau_{r,3} = (z - z_m) + (r_0 - r_{0,m}) \quad (6)$$

¹ Although the additional ground telescopes give no direct information about the observable, they may provide additional constraints that enable improved calibration of a variety of systematic errors in an observation.

Now, there are three equations in the four quantities x , y , z , and r_0 . However, a fourth useful equation expresses r in terms of its components

$$r = [x^2 + y^2 + z^2]^{1/2} \quad (7)$$

If one assumes that the clock initialization is done at the same time as the VLBI observations, $r_0 = r$ and $r_{0,m} = r_m$ in Eqs. (4-6). Then, the delay observables from the observations of the three radio sources carry sufficient information to determine the spacecraft position.

After combining Eqs. (4-7) and doing some algebra, one finds the following result for the spacecraft position:

$$\begin{aligned} x &= \frac{1}{2} [x_m - y_m - z_m - r_m + c(\tau_{r,1} - \tau_{r,2} - \tau_{r,3})] + \sqrt{\frac{K}{2}} \\ y &= \frac{1}{2} [y_m - x_m - z_m - r_m + c(\tau_{r,2} - \tau_{r,1} - \tau_{r,3})] + \sqrt{\frac{K}{2}} \\ z &= \frac{1}{2} [z_m - x_m - y_m - r_m + c(\tau_{r,3} - \tau_{r,1} - \tau_{r,2})] + \sqrt{\frac{K}{2}} \end{aligned} \quad (8)$$

Here, K is given by

$$\begin{aligned} K &= \frac{1}{2} [(x_m + y_m + z_m + r_m)^2 - c^2(\tau_{r,1}^2 + \tau_{r,2}^2 + \tau_{r,3}^2)] \\ &\quad + c\tau_{r,1}(c\tau_{r,2} + y_m + z_m - x_m + r_m) \\ &\quad + c\tau_{r,2}(c\tau_{r,3} + x_m + z_m - y_m + r_m) \\ &\quad + c\tau_{r,3}(c\tau_{r,1} + x_m + y_m - z_m + r_m) \end{aligned} \quad (9)$$

Equations (8) and (9) provide the means of finding the spacecraft position \vec{r} based on the model position \vec{r}_m and the three measured quantities $\tau_{r,1}$, $\tau_{r,2}$, and $\tau_{r,3}$.

The uncertainty in a component of the spacecraft position can be computed using Eqs. (8) and (9). Taking the x -component as an example, the uncertainty σ_x in x can be found from

$$\sigma_x^2 = \left(\frac{\partial x}{\partial \tau_{r,1}}\right)^2 \sigma_{\tau_{r,1}}^2 + \left(\frac{\partial x}{\partial \tau_{r,2}}\right)^2 \sigma_{\tau_{r,2}}^2 + \left(\frac{\partial x}{\partial \tau_{r,3}}\right)^2 \sigma_{\tau_{r,3}}^2 \quad (10)$$

Here, the uncertainty in the measurement of the delay residual for the measurement of radio source i is $\sigma_{\tau_{r,i}}$, and the uncertainties in the three radio-source measurements have been taken to be uncorrelated. After doing some algebra, the uncertainties in the components of the spacecraft position are found to be

$$\begin{aligned} \sigma_x^2 &= c^2 \left\{ \left[1 + \frac{x^2}{2K} - x\sqrt{\frac{2}{K}} \right] \sigma_{\tau_{r,1}}^2 \right. \\ &\quad \left. + \frac{y^2}{2K} \sigma_{\tau_{r,2}}^2 + \frac{z^2}{2K} \sigma_{\tau_{r,3}}^2 \right\} \\ \sigma_y^2 &= c^2 \left\{ \left[1 + \frac{y^2}{2K} - y\sqrt{\frac{2}{K}} \right] \sigma_{\tau_{r,2}}^2 \right. \\ &\quad \left. + \frac{x^2}{2K} \sigma_{\tau_{r,1}}^2 + \frac{z^2}{2K} \sigma_{\tau_{r,3}}^2 \right\} \\ \sigma_z^2 &= c^2 \left\{ \left[1 + \frac{z^2}{2K} - z\sqrt{\frac{2}{K}} \right] \sigma_{\tau_{r,3}}^2 \right. \\ &\quad \left. + \frac{x^2}{2K} \sigma_{\tau_{r,1}}^2 + \frac{y^2}{2K} \sigma_{\tau_{r,2}}^2 \right\} \end{aligned} \quad (11)$$

If the simplifying assumption is made that $\sigma_{\tau_{r,1}} = \sigma_{\tau_{r,2}} = \sigma_{\tau_{r,3}} \equiv \sigma_\tau$, Eq. (11) reduces to

$$\left. \begin{aligned} \sigma_x^2 &= \left[1 + \frac{r^2}{2K} - x\sqrt{\frac{2}{K}} \right] c^2 \sigma_\tau^2 \\ \sigma_y^2 &= \left[1 + \frac{r^2}{2K} - y\sqrt{\frac{2}{K}} \right] c^2 \sigma_\tau^2 \\ \sigma_z^2 &= \left[1 + \frac{r^2}{2K} - z\sqrt{\frac{2}{K}} \right] c^2 \sigma_\tau^2 \end{aligned} \right\} \quad (12)$$

The expression for K [Eq. (9)] can be given in terms of the actual spacecraft position components:

$$K = \frac{1}{2}(x + y + z + r)^2 \quad (13)$$

Combining Eqs. (12) and (13) yields the following result:

$$\sigma_x^2 = c^2 \sigma_\tau^2 \left[1 + \frac{r^2}{(x + y + z + r)^2} - \frac{2x}{x + y + z + r} \right] \quad (14)$$

For example, take the case where the spacecraft position is along the line of sight to one of the radio sources, say $x = r$ and $y = z = 0$. Then, evaluation of Eq. (14) shows that $\sigma_x = 0.5c\sigma_\tau$. Since the position error along the line of sight to the radio source and the position error along the line of sight to the tracking station are one and the same quantity in this case, the signature in the delay observable is doubled for a given position offset, or the position error is halved for a given delay measurement error. In order to show the range of values for σ_x , Fig. 2 is a plot of σ_x (in units of $c\sigma_\tau$) as a function of x/r , assuming that $z = 0$. Fig. 2(a) shows the results for $y > 0$, while Fig. 2(b) displays results for $y < 0$. Figures 3(a) and (b) show similar plots for σ_x under the assumption that $z = 0.5r$. In general, the results for the uncertainties along each of the three coordinate axes are similar for a given set of assumptions.

As Fig. 2 shows, σ_x diverges if $x = -r$ or $y = -r$ (or, as not shown in the figure, if $z = -r$), i.e., if the line of sight from the tracking station to the spacecraft is opposite to the direction to one of the natural radio sources. Consider the case where $x \approx -r$. In fact, $x = -r$ is not possible, since this would require the spacecraft radio telescope to look through the Earth to see the radio source hypothesized to be in the direction of the positive x -axis. However, it is possible that $x \approx -r$ in the case where the spacecraft is at a low elevation angle as seen from the tracking station and is several Earth radii distant. For $x = -r$, any error in the spacecraft position along the x -direction is exactly compensated for by an error in the time associated with the VLBI data, and an infinite position error would give no signature in the delay measurement.

The denominators of two terms in Eq. (14) vanish when $x + y + z = -r$, which includes (but is not limited to) the situation where the spacecraft position is opposite to the direction to one of the three radio sources. This is the equation of a plane. The spacecraft position must lie on the spherical surface at a range r from the tracking station; the intersection of the plane and that spherical surface is a circle on which the spacecraft position error becomes infinite. Therefore, if the hypothetical VLBI delay measurements were to be used to improve the spacecraft orbit

determination, they should be made at a time when the spacecraft position lies far from that circle of singularity. If there are two or more tracking stations available, it should be possible to select the one for which the position errors dictated by the geometry of the VLBI delay measurements are minimized and to use that station for the time transfer to the spacecraft.

IV. Estimate of VLBI Delay Measurement Errors

A number of practical limitations may prevent space VLBI observations from being as useful for navigation as indicated for the idealistic case treated above. This section provides a discussion of some of those limitations. For a number of effects, contributions to the delay error σ_τ (expressed in distance units, i.e., multiplied by c) are estimated independently, then combined in quadrature. It is important to recognize that the actual method of orbit determination would involve a multiparameter fit to spacecraft and radio-source positions, as well as other quantities such as troposphere, ionosphere, and clock parameters. Therefore, discussion of individual uncertainties as though they were completely separate from one another is an oversimplification but does serve to indicate the general limitations imposed by a variety of effects.

A. Limited Precision of Delay Measurements

The precision of the single-band delay measurements can be derived from the signal-to-noise ratio (SNR) of the VLBI observations using

$$\sigma_\tau = \frac{1}{2\pi\Delta\nu(\text{SNR})} \quad (15)$$

where $\Delta\nu$ is the observing bandwidth. Table 1 summarizes assumptions and predicted measurement precision for three different observing frequencies (1.6, 5, and 22 GHz) that will be used in the first generation of VLBI satellites. Assumptions are those for the current best guesses at the performance of the radio telescope on board Radioastron; the VSOP performance may be somewhat poorer. In distance units, the estimated delay precisions are 1.8 cm, 1.9 cm, and 4.5 cm at 1.6, 5, and 22 GHz, respectively. Radio sources with correlated flux densities as high as the assumed value of 0.5 Jy on baselines in the 40,000–80,000 km range probably will be rare or nonexistent at all three frequencies, so the above precisions will not be achievable on the longest baselines for most sources.

B. Ionospheric Propagation Errors

In the DSN, typical errors of 2–4 cm currently are achieved in the calibration of the ionospheric propagation delay of 8.4-GHz radio signals at the zenith [17], although there is some hope for improvement using signals from GPS satellites [18]. Even if it is assumed that the radio-source signal propagating to the spacecraft suffers no charged-particle delay because of the interplanetary medium, the delay measurement will still be corrupted by propagation through the Earth's ionosphere to the ground radio telescope. Scaling from the 8.4-GHz estimates (ionospheric delays are proportional to the inverse square of the frequency), the current ionospheric calibration errors will give respective zenith delay errors of 55–110 cm, 5–10 cm, and 0.3–0.6 cm at the three frequencies. (Radioastron also will operate at 300 MHz, where the ionospheric effects will be even larger.) At elevation angles of 30 deg, which will be much more common when several radio sources are observed in very different directions, the above delay errors should be doubled. The delay errors could be reduced substantially if there were a capability for simultaneous dual-frequency observations, as there is for ground-based astrometric and geodetic experiments. Such a capability does not exist on VSOP but might be available (at 1.6 and 5 GHz) on Radioastron. However, ground radio telescopes currently do not have dual-frequency capabilities at the space VLBI frequencies, implying that ionospheric propagation errors cannot be reduced by means of dual-frequency observations.

C. Tropospheric Propagation Errors

Troposphere fluctuations and errors in the static troposphere also will have a significant effect on delay measurements. Typical errors in the calibration of the zenith troposphere delay are about 4 cm, corresponding to 8 cm at a 30-deg elevation. By the mid-1990s, GPS calibrations have the potential for reducing these errors by a factor of 2–4 at ground radio telescopes, provided that GPS receivers are present at the telescopes. The error caused by troposphere fluctuations is on the order of 3 cm at low elevations, and not readily reducible using GPS data. In the future, this error might be reduced by using advanced water-vapor radiometers.

D. Earth-Orientation and Timing Errors

Errors in prediction of Earth orientation and Universal Time typically give errors of tens of centimeters or more in effective locations of tracking stations and radio telescopes on the Earth. However, a delay of several weeks will occur between the observations and the data correlation. With that delay, the combination of VLBI and GPS calibration measurements allows reconstruction of the Earth

orientation and timing parameters to better than 2 cm per component [19], implying similar errors in the delay measurements.

E. Radio-Source Position Errors

In the idealized case, it was assumed implicitly that the positions of the radio sources observed for orbit determination were known perfectly. In fact, they are not. For a priori position errors of 1 nrad that will be characteristic of the strongest compact radio sources in the mid-1990s, the delay measurement error will be 4 cm on a 40,000-km baseline and 8 cm on an 80,000-km baseline.

F. Summary of Delay Errors

Table 2 summarizes the minimum expected error contributions (in length units) to space-ground VLBI delay measurements for Radioastron, assuming a 40,000-km baseline. (The value of 40,000 km is used because for observations of three radio sources in orthogonal directions, it is not possible for the projected baselines in all three directions to be near 80,000 km.) This table assumes the sensitivity and the improved troposphere and Earth-orientation calibrations given in the above subsections. On longer baselines, the delay error contributed by the source position uncertainties will be larger, while the correlated flux densities and consequent sensitivities probably will be lower than assumed above. Even assuming no instrumental errors on the ground, the minimum delay errors derived from the rss of the individual error contributions are 110 cm at 1.6 GHz, 12 cm at 5 GHz, and 8 cm at 22 GHz. Inspection of Figs. 2 and 3 shows that the one-dimensional spacecraft position error derived from the VLBI data for the idealized radio-source observation strategy can range from 0.5 to more than 10 times the delay measurement error, depending on the exact geometry. These results are for the hypothetical case of simultaneous observations of three radio sources, and do not include the additional considerations discussed below in Section V.

V. Other Practical Considerations in Using VLBI Delay for Orbit Determination

A. Geometry of Radio-Source Observations

The ideal case considered above assumes observations of radio sources in three orthogonal directions. This simplification makes the analytical development tractable but probably is not necessary for improved orbit-determination results. In order to solve for the spacecraft position, it is likely that the three radio-source directions need only be linearly independent (i.e., not coplanar). However, it will be necessary to determine whether, in this more general

case, there are geometries in which the position error diverges as it does in the idealized case. If the assumption of orthogonal radio-source directions could be relaxed, this would help in several ways. First, it would give a much larger set of candidate sources for observation, which is important in view of the correlated-flux limitations mentioned previously. Second, it would make it more likely that a set of radio sources providing a reasonable geometry actually can be observed; spacecraft constraints will make it very difficult to observe sources in three orthogonal directions, none of which is near the line of sight to the tracking station. Third, the total amount of slewing necessary for the space radio telescope to observe the three radio sources would be reduced, a consideration whose importance is described further in the next subsection.

B. Nonsimultaneous Observations and Propagation to a Common Reference Time

The idealized case considered in this article includes the assumption that observations of three radio sources in very different directions could be made simultaneously. Of course, this assumption is completely unrealistic. For VSOP, the minimum time necessary for the space telescope to make a 90-deg slew to a new source, settle, and begin observations, is likely to be at least 60 minutes. Radioastron should slew much more rapidly and may be able to change sources in 15 minutes, so it is used here to derive the more optimistic result. Assuming 5-min integration times, three observations would take a total of 45 minutes, with reference times (scan midpoints) spanning 40 minutes. If the reference time were chosen to be the time of the second observation, the delay residuals at the times of the first and third observations must be propagated to the time of the second observation in order to solve for the spacecraft position at that time. Typical velocity uncertainties in the reconstructed spacecraft orbit will be about 1 cm/sec [11]. In the unlikely event that the velocity errors from one second to the next were completely uncorrelated, there would be an additional position error of at least 35 cm due to this error propagation over 20 minutes. If the correlation time for velocity errors in the reconstructed orbit is much longer than 1 sec, the position error at the reference time will be considerably larger than 35 cm. For example, if the correlation time for the velocity errors were 100 seconds or more, as seems likely, the position error due solely to the propagation to a reference time 20 minutes away would be at least 3 meters.

If the orbit at the reference time is propagated to the time of another VLBI observation that might be used for astrometric purposes, the spacecraft position error at that time will be still larger than the error at the reference time.

Uncertain spacecraft accelerations caused by mismodeling of the Earth's gravitational field, by errors in the model of solar pressure effects on the spacecraft, and by spacecraft maneuvers will serve to increase the spacecraft position error at times later (or earlier) than the epoch of the observations actually used for orbit determination. In addition, the above discussion is oversimplified because it assumes a fixed velocity accuracy that is used to propagate the position results found from the VLBI data. In reality, the VLBI data would be acquired simultaneously with two-way Doppler data, and both would be used to determine the spacecraft orbit. Detailed investigation of the orbit accuracy achievable in that case is beyond the scope of this article.

C. Tracking Continuity

In order for the three hypothetical VLBI observations to provide useful data for spacecraft trajectory determination, they must be referenced to the same clock. This implies, first, that the same tracking station must be used during all the VLBI observations. Second, it suggests that it may be necessary to maintain a continuous link to the tracking station between the VLBI observations. Without that continuity, there also will be clock breaks that will degrade the accuracy of the orbit determination. For a spacecraft in a fairly low orbit, such as VSOP, there will be a limited view period from a particular tracking station. In many instances, there will not be time to make three observations while the spacecraft is in view of the same station, particularly since long slews with VSOP may take up to 4 hours.

D. Possible Importance of Ranging

It was shown in Section III that in the absence of other data, three VLBI delay measurements are needed to improve on the knowledge of the spacecraft position. However, inspection of Eqs. (4-6) shows that a single VLBI delay observation can determine one component of the spacecraft position if the range r is determined accurately. Thus, a capability for ranging from the tracking stations to the VLBI spacecraft could be quite useful. In fact, any combination of three simultaneous ranging and VLBI delay measurements sampling three linearly independent directions should suffice to provide an accurate instantaneous position for the spacecraft. There is no fundamental reason that ranging observations from two different tracking stations could not be made at the same time that the spacecraft is making a VLBI observation of a radio source in a third direction. The two ranging measurements would provide accurate position components in two directions as well as supply an absolute clock time to the spacecraft. The accuracy of the VLBI delay observable still would be

limited by the effects summarized in Section IV above and by the clock error imposed by the limited accuracy of the ranging system. Two simultaneous ranging measurements with errors of 15 cm, when combined with the minimum error of 8 cm possible for a 22-GHz VLBI delay measurement, would give a total delay accuracy of about 23 cm, implying a spacecraft position uncertainty at least as large.

There are practical limitations and implications for the current generation of space VLBI satellites. VSOP and Radioastron will have limited ranging capability only from Japanese and Russian tracking stations, respectively. Radioastron will not be able to maintain a phase link during ranging observations, so simultaneous ranging/VLBI observations are not possible. Furthermore, neither spacecraft has multiple, independent downlink antennas. Therefore, it will not be possible to do simultaneous ranging from two different tracking stations. Reorientation of the downlink antenna would be necessary to do ranging from two different tracking stations; that reorientation would cause a break in the clock continuity.

E. Possible Benefits of Accurate Clocks On Board Space VLBI Satellites

Another possibility that would eliminate some of the difficulties in solving for the position of a space VLBI satellite would be the provision of a highly accurate clock on board the spacecraft. However, at centimeter wavelengths, monitoring of a two-way phase link from the ground, as was done in the first space VLBI experiments [3,4] and is planned for the first generation of dedicated missions, can provide a clock-rate accuracy nearly equivalent to that of the original clock on the ground. If a clock is flown on board the spacecraft, calibration of the absolute delay still will face the same difficulties as the calibration of the absolute delay for a clock transmitted from the ground. In either case, highly accurate ranging measurements are needed to fix the absolute clock time.

VI. Prospects (or Lack Thereof) for Improved Astrometry

The above analysis has shown that, with many caveats, it may be possible to achieve improvements in orbit determination through the use of space VLBI delay data. The degree of improvement that is achievable must be determined through numerical simulations, since the simple analytical model does not include many of the practical limitations discussed above. However, it is possible to address some aspects of the utility of the orbit determination improvement. In particular, a desirable consequence

of highly accurate orbit determination would be the ability to use the spacecraft as a platform for making astrometric VLBI observations more accurately than could be done using ground baselines alone. This section addresses the feasibility of that task.

It is critical to recognize that doing VLBI with a space radio telescope is a much more expensive and complicated task than ground-only VLBI experiments. Therefore, it makes no sense to do astrometry in space-ground VLBI experiments unless it can be shown that the accuracy will be significantly better than is possible for ground-only astrometry experiments. Expected improvement by a factor of 2 or more should be the minimum requirement for attempting an astrometric experiment using space VLBI.

The current delay precision achieved in the best ground-based astrometric VLBI experiments is approximately 15 psec [20], corresponding to a distance of 4.5 mm, a factor of about 20 better than the best that could be hoped for in space-ground VLBI using the first generation of dedicated VLBI satellites. The accuracy of the determination of VLBI baselines on the ground is approximately 1 cm in the best experiments [21], also a factor of about 35 better than the best that possibly could be expected for space VLBI in the 1990s. The delay precision achieved on the ground can be used to give much better baselines, in large part because of the ability to make many observations in different directions in a short period of time. This enables multiparameter fits that give excellent solutions for such quantities as clocks and atmospheric delays. Space VLBI observations for astrometric purposes will not be competitive with ground-only observations until they are capable of such flexibility.

The astrometric angular precision possible for wide-angle VLBI astrometry, σ_ϕ , can be estimated from the baseline error σ_B and the baseline length B as

$$\sigma_\phi \sim \frac{\sigma_B}{B} \quad (16)$$

In ground-only experiments, a baseline error of 1 cm over a 10,000-km baseline corresponds to an angular error of 1 nrad. For Radioastron, the minimum possible baseline error of 35 cm over a baseline length of 40,000 km gives an angular accuracy of 9 nrad in the best possible case, not competitive with current ground-only VLBI.

It also is important to consider the possibility that accurate differential astrometric VLBI from space might be done over narrow angles without knowledge of the spacecraft orbit to within a few centimeters. One could imagine

successive, quasi-simultaneous observations of two radio sources separated in the sky by the small angle ϕ . For ϕ much smaller than a radian, the differential astrometric error in the source positions would be given approximately by

$$\sigma_\phi \approx \phi \frac{\sigma_B}{B} \quad (17)$$

The extra factor of ϕ compared with Eq. (16) is due to the fact that the effect of the baseline error becomes smaller and smaller for decreasing angular separations of the radio sources, as a larger fraction of the delay error cancels when the delay observations of the two radio sources are differenced. For angular separations on the order of 0.5 deg, differential astrometry at the 0.1-nrad level has been achieved using ground radio telescopes [22]. The space-ground VLBI differential astrometry on a 40,000-km baseline could yield a factor-of-two improvement at this angular separation only if the spacecraft position were known with an accuracy of ~ 20 cm. This accuracy does not seem achievable using space VLBI delay data.

VLBI observations of two sources in the same beams of the radio telescopes might provide further error reduction and the ability to do highly accurate differential VLBI. Given the sensitivity limitations of the first generation of space VLBI telescopes, there will be few (or no) pairs of continuum sources with separations well under a degree that can be observed. Thus the narrow-angle differential VLBI might be limited to spectral-line studies of the separation of spots in water masers. Separations of different water-maser complexes in nearby external galaxies might be studied; for same-beam VLBI involving a 70-m ground telescope, these complexes would need to be separated by less than 0.01 deg. Individual water-maser sources typically span regions of 100 nrad (6×10^{-6} deg) in this galaxy to 1 nrad or less in external galaxies, so measurements of spot separations within these individual complexes also might be possible.

Ten years ago, ground-based observations of the quasar pair 1038+528A and B, separated by about 0.01 deg, achieved differential astrometric accuracy of 0.02 nrad [23]. For a 40,000-km baseline and a 0.01-deg separation, Eq. (17) predicts that the space-ground VLBI astrometry could achieve this accuracy if the spacecraft position error were less than about 5 meters. Thus, it is conceivable that the space observations might compete at 0.01-deg source separations, although the current potential for ground-only results may be significantly better than the accuracy obtained in the early 1980s. However, lack of frequency tunability will prevent Radioastron from observing

most extragalactic water masers, and the shorter baselines sampled using VSOP will require position errors less than about two meters for such astrometric measurements to be more useful than the ground VLBI observations.

For source separations of 500 nrad (3×10^{-5} deg) within individual water maser complexes, Eq. (17) predicts that orbit determination accuracies of 40 m on a 20,000-km baseline would provide the potential for differential astrometric accuracies better than 10^{-3} nrad. Such orbital accuracies should be achievable using standard two-way Doppler tracking [11]. Furthermore, the sensitivities of VSOP and Radioastron probably would limit the possible astrometric accuracy to $\sim 5 \times 10^{-3}$ nrad for water masers [24]. Thus, improved orbit determination is not needed for accurate differential astrometry within individual maser complexes.

VII. Summary

An analysis has been performed to determine the information content of VLBI delay measurements on a space-ground baseline for improved orbit determination of an Earth-orbiting VLBI satellite. In the idealized case of simultaneous observations of three different radio sources in mutually orthogonal directions, expressions have been given for the expected error in the orbit determination from the VLBI data alone. Given a number of random and systematic error sources that will be very difficult to reduce in space VLBI, the delay measurement precision even in the idealized case will be at least 8 cm for the VLBI spacecraft scheduled for launch in the mid-1990s; this often corresponds to a position error much larger than 8 cm. Other practical aspects of space VLBI, particularly the long time required between observations, will make the orbit determination less accurate than the prediction for the ideal case. The long time between observations would give rise to a minimum position error of at least 35 cm, with an error of at least a few meters far more likely. Combination of VLBI data with simultaneous highly accurate ranging data from two different tracking stations shows promise for an improved determination of the instantaneous spacecraft position for a second generation of dedicated space VLBI observatories.

It is unrealistic to expect that an instantaneous position accuracy as good as 35 cm can be derived from the space VLBI delay data acquired by the first generation of space VLBI satellites; the achievable accuracy probably will be no better than 3 meters. Even under the optimistic assumption of 35-cm position accuracy, space-ground VLBI

cannot compete with ground-based VLBI for astrometry over angles of 0.1 deg or larger. It is conceivable that differential astrometric observations over smaller separations might be useful, although the lack of sensitivity of the first generation of space VLBI telescopes probably limits this utility to relative measurements of water-maser spots. Within individual maser complexes, the accuracy of the spacecraft position usually is not the limiting factor for astrometric accuracy, and orbit determination using two-way Doppler data may suffice.

Several changes could be implemented for the second generation of space VLBI observations that would provide enhanced capabilities for improved orbit determination using the space-ground VLBI delay data. First, the capability for simultaneous dual-frequency operations is needed for both the spacecraft radio telescope and the large ground telescopes, in order to solve for and eliminate the charged particle effects on signal propagation to the ground telescopes. (A number of ground telescopes currently can make simultaneous observations at 2.3 and 8.4 GHz, but this capability is of no use for space VLBI satellites that do not observe at those frequencies.) Second, the most accurate possible calibrations must be used to minimize the errors in knowledge of Earth orientation and tropospheric delay. Third, the space radio telescope must be more sensitive, implying some combination of

larger bandwidth, larger diameter, and lower system temperature. Fourth, the space radio telescope must have a capability for very rapid slewing and for observing a large number of sources in a relatively short period of time. Fifth, a highly accurate ranging system, including the capabilities for simultaneous ranging from more than one direction and for simultaneous ranging and VLBI data acquisition, is necessary to realize the improvements discussed in this article. This capability could be supplied either by ranging from the ground or by a satellite system, such as GPS, preferably with higher orbits. Sixth, solar-pressure and gravitational-field models would have to be made more accurate to improve the propagation of trajectory measurements to a specific reference time. Seventh, simultaneous observations with a large number of ground telescopes could be helpful in solving for systematic calibration errors.

With all the improvements listed above, it might be possible to use space VLBI delay data to determine an instantaneous spacecraft position with an accuracy of tens of centimeters. However, this still is not accurate enough for space-ground VLBI to compete with ground-only VLBI in making global astrometric measurements. Instead, space VLBI can be used far more productively for imaging and other astrophysical measurements that depend on the longer baselines available, but do not require highly accurate measurements of the absolute delay.

Acknowledgments

The author thanks Roger Linfield, Carl Christensen, Jeff Estefan, and Sam Thurman for very helpful comments on earlier versions of this article.

References

- [1] A. R. Thompson, J. M. Moran, and G. W. Swenson, Jr., *Interferometry and Synthesis in Radio Astronomy*, New York: John Wiley and Sons, pp. 247-313, 1986.
- [2] M. J. Reid and J. M. Moran, editors, *The Impact of VLBI on Astrophysics and Geophysics*, International Astronomical Union Symposium No. 129, Dordrecht, the Netherlands: Kluwer, 1988.
- [3] G. S. Levy, et al., "Very Long Baseline Interferometric Observations Made With an Orbiting Radio Telescope," *Science*, vol. 234, pp. 187-189, October 10, 1986.

- [4] G. S. Levy, et al., "VLBI Using a Telescope in Earth Orbit. I. The Observations," *Astrophysical Journal*, vol. 336, pp. 1098-1104, 1989.
- [5] R. P. Linfield, et al., "VLBI Using a Telescope in Earth Orbit. II. Brightness Temperatures Exceeding the Inverse Compton Limit," *Astrophysical Journal*, vol. 336, pp. 1105-1112, 1989.
- [6] R. P. Linfield, et al., "15 GHz Space VLBI Observations Using an Antenna on a TDRSS Satellite," *Astrophysical Journal*, vol. 358, pp. 350-358, 1990.
- [7] H. Hirabayashi, "VLBI Activities in Japan and a Projected Space-VLBI Program," in *The Impact of VLBI on Astrophysics and Geophysics*, International Astronomical Union Symposium No. 129, edited by M. J. Reid and J. M. Moran, Dordrecht, the Netherlands: Kluwer, pp. 449-456, 1988.
- [8] N. S. Kardashev and V. I. Slysh, "The Radioastron Project," in *The Impact of VLBI on Astrophysics and Geophysics*, International Astronomical Union Symposium No. 129, edited by M. J. Reid and J. M. Moran, Dordrecht, the Netherlands: Kluwer, pp. 433-440, 1988.
- [9] *IVS: An Orbiting Radio Telescope*, Report on the Assessment Study, European Space Agency, January 1991.
- [10] G. Tang, "A Short Note on the High-Precision Navigation of Quasat," in *Proceedings of Workshop on Quasat*, Gros Enzersdorf, Austria, ESA SP-213, pp. 185-186, June 18-22, 1984.
- [11] J. A. Estefan, "Precise Orbit Determination of High-Earth Elliptical Orbiters Using Differenced Doppler and Range Measurements," *TDA Progress Report 42-106*, vol. April-June 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 1-22, August 15, 1991.
- [12] S. M. Lichten and J. A. Estefan, "High Precision Orbit Determination for High-Earth Elliptical Orbiters Using the Global Positioning System," paper AIAA 90-2954, AIAA/AAS Astrodynamics Conference, Portland, Oregon, August 20-22, 1990.
- [13] I. Fejes, I. Almár, J. Ádám, and Sz. Mihály, "On Astrometric and Geodynamic Aspects of Space VLBI," paper presented at the Intercosmos Scientific Conference, Szentendre, Hungary, May 1987.
- [14] I. Fejes and Sz. Mihály, "Application of Space-VLBI to Satellite Dynamics," *Advances in Space Research*, vol. 11, no. 2, pp. 429-437, 1991.
- [15] C. S. Christensen and J. A. Estefan, "NASA Orbit Determination Capability," in *Frontiers of VLBI*, Proceedings of the International VSOP Symposium, edited by H. Hirabayashi, M. Inoue, and H. Kobayashi, Tokyo: Universal Academy Press, pp. 119-123, 1991.
- [16] L. R. D'Addario, "Time Synchronization in Orbiting VLBI," *IEEE Transactions on Instrumentation and Measurement*, vol. 40, no. 3, pp. 584-590, June 1991.
- [17] H. N. Royden, R. B. Miller, and I. A. Buennagel, "Comparison of NAVSTAR Satellite L-Band Ionospheric Calibrations with Faraday Rotation Measurements," *Radio Science*, vol. 19, no. 3, pp. 798-804, May-June 1984.
- [18] G. E. Lanyi and T. Roth, "A Comparison of Mapped and Measured Total Ionospheric Electron Content Using Global Positioning System and Beacon Satellite Observations," *Radio Science*, vol. 23, pp. 483-492, July-August 1988.

- [19] U. J. Lindqwister, A. P. Freedman, and G. Blewitt, "A Demonstration of Centimeter-Level Monitoring of Polar Motion With the Global Positioning System," *TDA Progress Report 42-108*, vol. October-December 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 1-9, February 15, 1992.
- [20] J. R. Ray and B. E. Corey, "Current Precision of VLBI Multi-Band Delay Observables," *Proceedings of the AGU Chapman Conference on Geodetic VLBI: Monitoring Global Change*, Washington, DC, April 22-26, 1991.
- [21] T. A. Clark, W. G. Melbourne, C. Reigber, L. E. Young, and T. P. Yunck, "Instrumentation: Microwave Techniques," *Proceedings of an International Workshop held at Erice, The Interdisciplinary Role of Space Geodesy, Lecture Notes in Earth Sciences* (edited by I. I. Mueller and S. Zerbini), vol. 22, Berlin, Germany: Springer-Verlag, pp. 148-162, 1989.
- [22] N. Bartel, M. I. Ratner, I. I. Shapiro, T. A. Herring, and B. E. Corey, "Proper Motion of Components of the Quasar 3C 345," in *VLBI and Compact Radio Sources*, International Astronomical Union Symposium No. 110, edited by R. Fanti, K. Kellermann, and G. Setti, Dordrecht, the Netherlands: Reidel, pp. 113-116, 1984.
- [23] J. M. Marcaide and I. I. Shapiro, "High Precision Astrometry via Very-Long-Baseline Radio Interferometry: Estimate of the Angular Separation Between the Quasars 1038+528A and B," *Astronomical Journal*, vol. 88, pp. 1133-1137, 1983.
- [24] B. F. Burke, et al., *US Participation in VSOP and Radioastron*, Report of the U.S. OVLBI Science Consulting Group, pp. 21-22, August 1, 1989.

Table 1. Space VLBI delay-measurement precision.

Observing frequency, GHz	1.6	5.0	22
Correlated flux density, Jy	0.5	0.5	0.5
Observing bandwidth, MHz	32	32	32
Integration time, sec	300	300	300
Ground telescope			
Diameter, m	70	70	70
Efficiency	0.60	0.60	0.50
System temperature, K	35	35	50
Space telescope			
Diameter, m	10	10	10
Efficiency	0.50	0.50	0.30
System temperature, K	60	70	135
Sensitivity			
σ_{τ} , nsec	0.06	0.06	0.15
$c\sigma_{\tau}$, cm	1.8	1.9	4.5

Table 2. Components of space VLBI delay-measurement errors (cm) on a 40,000-km baseline.

Observing frequency, GHz	1.6	5.0	22
System noise	1.8	1.9	4.5
Ionosphere	110-220	10-20	0.6-1.2
Static troposphere	2-4	2-4	2-4
Fluctuating troposphere	3	3	3
Earth orientation	3	3	3
Radio-source position	4	4	4
RSS	110-220	12-21	8

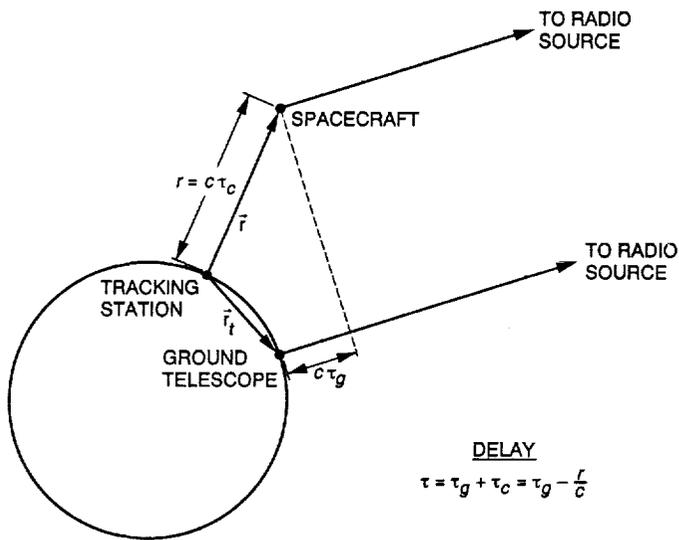


Fig. 1. Geometry of the delay measurement on a space-ground VLBI baseline.

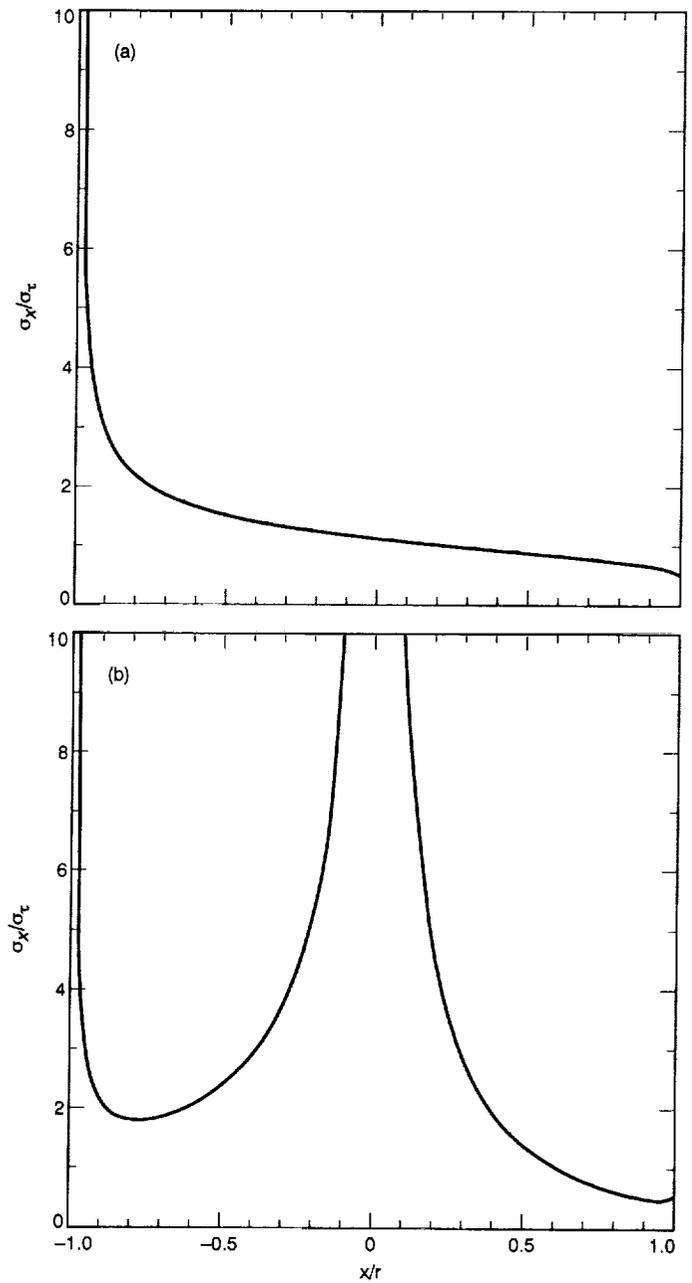


Fig. 2. Predicted minimum (one-dimensional) spacecraft position error, in units of the delay measurement error, for a highly idealized set of space-ground VLBI delay measurements: (a) $y > 0$ and $z = 0$ and (b) $y < 0$ and $z = 0$.

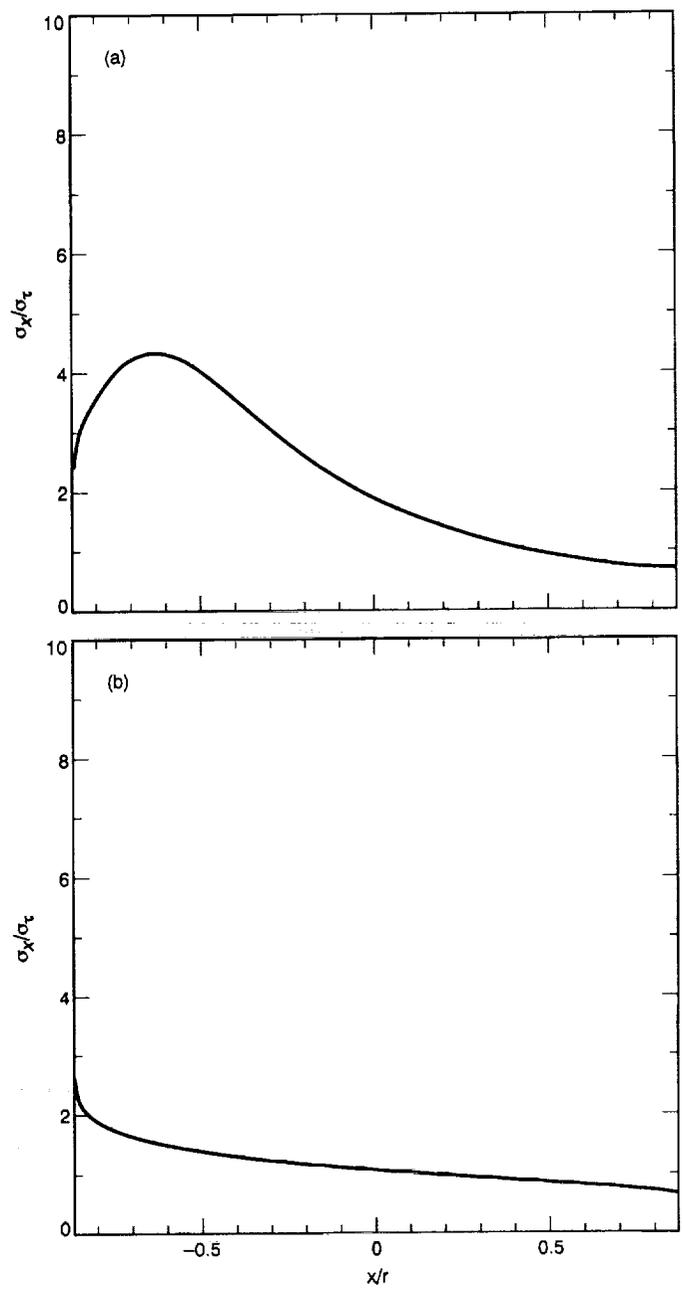


Fig. 3. Predicted minimum (one-dimensional) spacecraft position error, in units of the delay measurement error, for a highly idealized set of space-ground VLBI delay measurements: (a) $y > 0$ and $z = 0.5$ and (b) $y < 0$ and $z = 0.5$.

57-04

128 N 93-19420

p. 12

Modeling the Global Positioning System Signal Propagation Through the Ionosphere

S. Bassiri

Sharif University of Technology, Tehran, Iran

G. A. Hajj

Tracking Systems and Applications Section

Based on realistic modeling of the electron density of the ionosphere and using a dipole moment approximation for the Earth's magnetic field, one is able to estimate the effect of the ionosphere on the Global Positioning System (GPS) signal for a ground user. The lowest order effect, which is on the order of 0.1-100 m of group delay, is subtracted out by forming a linear combination of the dual frequencies of the GPS signal. One is left with second- and third-order effects that are estimated typically to be ~0-2 cm and ~0-2 mm at zenith, respectively, depending on the geographical location, the time of day, the time of year, the solar cycle, and the relative geometry of the magnetic field and the line of sight. Given the total electron content along a line of sight, the authors derive an approximation to the second-order term which is accurate to ~90 percent within the magnetic dipole moment model; this approximation can be used to reduce the second-order term to the millimeter level, thus potentially improving precise positioning in space and on the ground. The induced group delay, or phase advance, due to second- and third-order effects is examined for two ground receivers located at equatorial and mid-latitude regions tracking several GPS satellites.

I. Introduction

The Global Positioning System (GPS) consists of 24 satellites, evenly distributed in 6 orbital planes around the globe, at an altitude of about 20,200 km. Precise positioning of the GPS satellites, as well as ground and space users, is now reaching a few parts in 10^9 [1-6]. In addition, the GPS has been heavily utilized in a host of geodetic and other applications. These include seismic tectonic motions [7-9], Earth orientation studies [10,11], gravimetry [12], at-

mospheric water vapor calibration [13,14], and ionospheric monitoring [15]. Precise positioning and other GPS-based applications, however, require a very good understanding of all effects on the GPS signal as it propagates through the Earth's atmosphere, so that all effects can be solved for or modeled.

The GPS transmits two right-hand circularly polarized (RCP) signals at L-band frequencies: L1 at 1574.42 MHz

and L2 at 1227.6 MHz, which correspond to wavelengths of 19.0 cm and 24.4 cm, respectively. These are modulated by a pseudorandom precision code (P-code) at a frequency of 10.23 MHz [16]. [The additional lower frequency course acquisition (C/A) modulation is not of concern here.] A single measurement for a given transmitter and receiver pair will consist of four observables that will be denoted here by L_1 , L_2 for the accumulated carrier phase measurements at the two frequencies and P_1 , P_2 for the corresponding P-code pseudorange. In addition to the geometric range delay, the signals will experience delays, or phase advances, due to the presence of the ionosphere and neutral atmosphere.

The delay due to the neutral atmosphere is the same for all observables; its effect is on the order of 2 m and can be solved for to better than a centimeter [13,14]. However, due to the dispersive nature of the ionosphere, the group delay caused by it (or phase advance) is frequency dependent, and is on the order of 0.1–100 m, depending on the time of day, the time of year, and the solar cycle. If the ionospheric effect on signal delay (or advance) is expanded in powers of inverse frequencies, then the lowest order term ($1/f^2$), by far the most dominant, can be solved for and subtracted out by virtue of the dual frequencies of the GPS. Remaining higher order terms are on the order of submillimeters to several centimeters, which remain embedded in the signal and contribute to range and accumulated phase errors. While the first-order term depends simply on the total electron content (TEC), namely the integrated electron density inside a columnar cylinder of unit area between the transmitter and the receiver, higher order terms depend on the coupling between the Earth's magnetic field and the electron density everywhere along the line of sight. In order to estimate the higher order effects on the GPS observables, the authors modeled the ionosphere by a sum of Chapman layers and the Earth's magnetic field by that of a dipole moment. Such a model will make it possible to estimate higher order ionospheric effects at different geographical locations on the ground as well as their sensitivity to the electron density distribution. It will be demonstrated that knowledge of the TEC can be used to calibrate most of the second-order effect and reduce P-code and phase measurement errors to a few millimeters.

Due to the inhomogeneity of the propagation medium, the GPS signal does not travel along a perfectly straight line. Moreover, since the medium is dispersive, the two frequencies will take two slightly different paths. By applying the empirical formula given by Brunner and Gu [17] on the ionospheric model used below, the residual range error between the dual-frequency corrected range and the

true range, due to bending alone, is estimated to be ~ 4 mm at a 10-deg elevation angle and less than a millimeter for elevations above 30 deg. The bending effect will be ignored in the following analysis; the two signals will be assumed to travel along the same straight line.

A more elaborate modeling of higher order ionospheric effects, where bending is taken into account, has been considered by Brunner and Gu [17]; see also [18]. In their paper, the international geomagnetic reference fields (IGRF) and a Chapman profile of the ionosphere were used to estimate the residual range error. They also proposed an improved linear combination that corrects for the second- and third-order terms, as well as for bending. Their improved linear combination requires knowledge of N_m and h_m , the electron density peak and its altitude, respectively. In this article, the second- and third-order terms are considered separately. Here the authors estimate that the second-order term is dominant most of the time over the third-order and the curvature terms. A method of modeling the second-order effect based on a thin shell model of the ionosphere and a dipole magnetic field is suggested. The second- and third-order errors are considered at different geographical locations while tracking different satellites. It is demonstrated that knowledge of the TEC alone can be used to reduce the higher order effects to a few millimeters.

II. Earth's Ionosphere

The Earth's ionosphere extends from an altitude of about 80 to 1000 km. It is a macroscopically neutral ionized gas consisting principally of free electrons, ions, and neutral atoms or molecules. Ions in that region are 2000 to 60,000 times more massive than electrons. Thus, at the frequencies used for radio communication, the range of movement of an ion caused by the electric field of a radio wave is smaller than that of an electron by about the same factor. This implies that the ions can, for most purposes, be ignored [19].

The electron density profile exhibits several distinct regions (E, F1, and F2) as a result of the competing processes of particle production, loss, and transport. The maximum electron densities (10^{12} to 10^{13} m^{-3}) are observed at the F2 peak; the peak altitude ranges from 250 to 350 km at mid-latitudes and from 350 to 500 km at equatorial latitudes. The F1 region, which is present during the day but absent at night, has a peak near the 200-km altitude and is 3–5 times smaller than that of F2. The E peak density is about one order of magnitude smaller than the F2

peak and is typically located at the 100- to 120-km altitude. During daytime there is also a D region below the E region, with a peak at the 80-km altitude [20].

III. Propagation of Electromagnetic Waves in the Ionosphere

When a magnetostatic field \vec{B}_0 is applied to a plasma, the plasma becomes anisotropic for the propagation of electromagnetic waves. That is, the scalar dielectric constant of the plasma is transformed into a tensor. To study the propagation and polarization properties of a plane monochromatic wave in a magnetically biased homogeneous lossless plasma, the plasma is regarded as a continuous medium whose conductivity is zero, whose permeability is equal to that of a vacuum, and whose dielectric constant is a tensor. By solving the Helmholtz wave equation subject to proper constitutive relations, one can obtain the expressions for the fields and for the index of refraction. The index of refraction, n , for the Earth's ionosphere is given by the Appleton-Hartree formula [21], as follows:

$$n_{\pm}^2 = 1 - \frac{2X(1-X)}{2(1-X) - Y_{\perp}^2 \pm \sqrt{Y_{\perp}^4 + 4(1-X)^2 Y_{\parallel}^2}} \quad (1)$$

where

$$X = \left(\frac{f_p}{f}\right)^2 = \frac{(Ne^2/4\pi^2\epsilon_0 m)}{f^2} \quad (2)$$

$$Y_{\perp} = Y \sin \theta_B; \quad Y_{\parallel} = Y \cos \theta_B \quad (3)$$

$$Y = \left(\frac{f_g}{f}\right) = \frac{(|e|B_0/2\pi m)}{f} \quad (4)$$

N is the number density of electrons; e and m are the electron charge and mass, respectively; ϵ_0 is the permittivity of the free space; f_p , f_g , and f are the plasma, gyro, and carrier frequencies, respectively; and θ_B is the angle between the Earth's magnetic field, \vec{B}_0 , and the direction of propagation of the wavefront, \vec{k} . By definition, $\vec{Y} = e\vec{B}_0/2\pi fm$, and since e is negative, \vec{Y} is antiparallel to \vec{B}_0 . The plasma frequency is the natural frequency of oscillation for a slab of neutral plasma with density N after the electrons have been displaced from the ions and are allowed to move freely. The gyro frequency is the natural frequency at which free electrons circle around the magnetic field lines. For the Earth's ionosphere, with $N = 10^{12}$

electrons/m³, the plasma frequency $f_p \approx 8.9$ MHz. The gyro frequency for an electron in the Earth's magnetic field (2×10^{-5} tesla) is $f_g \approx 0.59$ MHz.

The plus and minus signs of Eq. (1) correspond to the ordinary and extraordinary wave modes of propagation, respectively. In general, these two waves are elliptically polarized with left and right senses of rotation, respectively. As a result of different phase velocities of the two waves, the total wave (the sum of ordinary and extraordinary waves) undergoes Faraday rotation as it passes through the ionosphere. When the carrier frequency is large, as compared with plasma and gyro frequencies, the principal modes of propagation are dominantly circularly polarized. This is the case for the GPS carrier frequencies.

Assuming that $Y \ll 2|\cos \theta_B|(1-X)/\sin^2 \theta_B$, the index of refraction can be expanded in inverse powers of frequency. For the GPS carrier frequencies, one has $(f_p/f) = 5.65 \times 10^{-3}$ and 7.25×10^{-3} , as well as $(f_g/f) = 3.75 \times 10^{-4}$ and 4.81×10^{-4} for L1 and L2, respectively. Therefore, the stated assumption is valid for GPS frequencies up to a value of $\theta_B \approx 89$ deg. The expansion of Eq. (1) up to the fourth inverse powers of frequency gives

$$n_{\pm} = 1 - \frac{1}{2}X \pm \frac{1}{2}XY|\cos \theta_B| - \frac{1}{4}X \left[\frac{1}{2}X + Y^2(1 + \cos^2 \theta_B) \right] \quad (5)$$

The second, third, and fourth terms on the right-hand side of Eq. (5) are proportional to the inverse square, inverse cube, and inverse quartic powers of frequency, respectively. The two values of n refer to the ordinary (+) and extraordinary (-) waves. At this point it should be noted from Eq. (5) that the index of refraction is smaller than unity, which corresponds to a phase velocity greater than the speed of light (phase advance). The group refractive index, on the other hand, given by $n^{\text{group}} = n + f(dn/df)$, can be written as

$$n_{\pm}^{\text{group}} = 1 + \frac{1}{2}X \mp XY|\cos \theta_B| + \frac{3}{4}X \left[\frac{1}{2}X + Y^2(1 + \cos^2 \theta_B) \right] \quad (6)$$

The group delay of a signal passing through the ionosphere, relative to vacuum as a reference, can be rewritten as

$$\tau_{\pm}^{\text{group}} = \frac{1}{c} \int (n_{\pm}^{\text{group}} \cos \alpha - 1) dl \quad (7)$$

where dl is an element of length along the line of sight, c is the velocity of light in a vacuum, and α is the angle between the wave normal and the ray direction. This angle has significance in anisotropic media, where the direction of the wave normal is, in general, different from the direction of energy propagation. Angle α can be found from the following relation: $\tan \alpha = (1/n)\partial n/\partial \theta_B$. By using Eq. (5) and the definition of α , it is easy to show that for the GPS carrier frequencies $\cos \alpha$ is essentially unity. By using Eqs. (5)–(7), the GPS observables can be written [ignoring the left-hand circularly polarized (LCP) component of the GPS signal, which has <0.35 percent and <2.5 percent of the total power, for L1 and L2, respectively] as

$$P_1 = \rho + \frac{q}{f_1^2} + \frac{s}{f_1^3} + \frac{r}{f_1^4} \quad (8a)$$

$$P_2 = \rho + \frac{q}{f_2^2} + \frac{s}{f_2^3} + \frac{r}{f_2^4} \quad (8b)$$

$$L_1 = \rho + n_1 \lambda_1 - \frac{q}{f_1^2} - \frac{1}{2} \frac{s}{f_1^3} - \frac{1}{3} \frac{r}{f_1^4} \quad (9a)$$

$$L_2 = \rho + n_2 \lambda_2 - \frac{q}{f_2^2} - \frac{1}{2} \frac{s}{f_2^3} - \frac{1}{3} \frac{r}{f_2^4} \quad (9b)$$

where

$$q = \frac{1}{2} \int f_p^2 dl = 40.3 \int N dl = 40.3 \text{ TEC} \quad (10)$$

$$s = \int f_g f_p^2 |\cos \theta_B| dl = 7527c \int N B_0 |\cos \theta_B| dl \quad (11)$$

$$r = 2437 \int N^2 dl + 4.74 \times 10^{22} \int N B_0^2 (1 + \cos^2 \theta_B) dl \quad (12)$$

TEC is the total electron content along the line of sight, and λ is the operating wavelength. In Eqs. (8) and (9), ρ corresponds to the geometrical distance plus all the nondispersive terms that are common to both frequencies, such as clocks, transmitter and receiver delays, and the neutral atmospheric delay. In Eq. (9), $n_1 \lambda_1$ and $n_2 \lambda_2$ correspond to unknown integer numbers of cycles that are constants for a given transmitter and receiver pair over a continuous tracking period. In addition to the terms shown on

the right-hand side of Eqs. (8) and (9), there are terms due to multipath, thermal noise, phase center variations, and a transmitter and receiver relative geometry dependent term; however, these are not the subject of this study, and are omitted from Eqs. (8) and (9).

IV. Ionospheric Layers and Geomagnetic Field Models

To proceed with the computation of the higher order delays, one has to assume models for the electron density, N , and the Earth's magnetic field, B_0 . For the electron density distribution, the Chapman layer model is chosen. This model is derived by assuming a homogeneous composition for air at a constant temperature. The curvature of the Earth is neglected, and it is assumed that the atmosphere is horizontally stratified and the scale height H_s is independent of height. As the solar radiation travels downward through the atmosphere, it is absorbed and hence ionization is produced. The rate of electron production is a function of height above mean sea level h and the sun's zenith angle χ , which is the angle between the ray from the sun and the zenith. From considerations of the production of electrons by photoionization and their removal by recombination, the following formula for the electron density distribution can be obtained [22]:

$$N = N_{\max} \exp \left[\frac{1}{2} (1 - z - e^{-z} \sec \chi) \right] \quad (13)$$

where N_{\max} is the maximum value of the electron density at an altitude of h_{\max} and $z = (h - h_{\max})/H_s$. When χ is near 90 deg, as near sunrise and sunset, the plane Earth approximation fails. To correct for this, $\sec \chi$ in Eq. (13) is replaced with the grazing incidence function $\text{Ch}(x, \chi)$. This function, which applies accurately only to a spherically symmetric atmosphere with H_s independent of height, can be expressed as

$$\begin{aligned} \text{Ch}(x, \chi) = & \left(\frac{1}{2} \pi x \sin \chi \right)^{1/2} e^{1/2x \cos^2 \chi} \\ & \times \left[1 \pm \text{erf} \left(\frac{1}{2} x \cos^2 \chi \right)^{1/2} \right] \quad (14) \end{aligned}$$

where $x = (R_E + h)/H_s$, R_E is the Earth's radius, and $\text{erf}(\cdot)$ is the error function. The plus (minus) sign refers

to $\chi > 90$ deg ($\chi < 90$ deg). Figure 1 is a plot of the electron density distribution versus height for two different solar zenith angles $\chi = 0$ deg and $\chi = 64$ deg. In obtaining this distribution, three different Chapman layers were added together so that the distribution can resemble the ionospheric F_2 , F_1 , and E layers: the E layer with a maximum at 110 km, the F1 layer with a maximum at 210 km, and the F2 layer with a maximum at 350 km. This figure is representative of a daytime profile typical of a year near sunspot maximum. The D layer, which is normally present during the daytime, is not included. During nighttime, the F1 layer disappears and the electron density for a given height is about 10–100 times smaller than that of daytime. In a solar minimum, the same features (D, E, F1, and F2 layers) are preserved with the electron density scaled down roughly by a factor of 10.

Next, one must model the Earth's magnetic field. A first approximation to the geomagnetic field near the surface of the Earth is an Earth-centered dipole with its axis tilted to intersect the Earth at 78.5 deg N latitude, 291.0 deg E longitude, which corresponds to the geomagnetic north pole; and at 78.5 deg S latitude, 111.0 deg E longitude, which corresponds to the geomagnetic south pole [20] (see Fig. 2).

At this point one must distinguish between two reference frames with a common origin at the Earth's center. The geodetic frame is Earth-fixed and is given by $\hat{x}, \hat{y}, \hat{z}$, where \hat{z} is along the Earth's spin axis, and \hat{x} is pointing toward 0 deg longitude. The geomagnetic frame, on the other hand, is obtained by first rotating the geodetic frame by an angle $\beta = 291$ deg around its \hat{z} axis, and then applying a second rotation by an angle $\delta = 11.5$ deg around the new \hat{y}_m axis (Fig. 3). This geomagnetic frame is denoted by $\hat{x}_m, \hat{y}_m, \hat{z}_m$ and is constructed so that \hat{z}_m is along the magnetic dipole. A vector transformation from the geodetic to the geomagnetic frame is given by

$$\vec{V}_m = \begin{pmatrix} \cos \delta \cos \beta & \cos \delta \sin \beta & -\sin \delta \\ -\sin \beta & \cos \beta & 0 \\ \sin \delta \cos \beta & \sin \delta \sin \beta & \cos \delta \end{pmatrix} \vec{V} \quad (15)$$

At a point on the Earth's surface, local geodetic east, north, and vertical are denoted by $\hat{X}, \hat{Y}, \hat{Z}$, and geomagnetic east, north, and vertical are denoted by $\hat{X}_m, \hat{Y}_m, \hat{Z}_m$ (Fig. 3). The magnetic field vector is given by

$$\vec{B}_0 = B_g \left(\frac{R_E}{r_m} \right)^3 \sin \theta_m \hat{Y}_m - 2B_g \left(\frac{R_E}{r_m} \right)^3 \cos \theta_m \hat{Z}_m \quad (16)$$

where r_m is the radial distance, and θ_m is the magnetic colatitude. The value B_g is the amplitude of the magnetic field at the Earth's surface at the magnetic equator, and is equal to 3.12×10^{-5} tesla.

V. Analysis

A. First-Order Effect

According to Eqs. (8)–(10), the first-order ionospheric delay can be written as $4.48 \times 10^{-16} \lambda^2 \text{TEC}$ (meters). For the GPS L1 and L2 frequencies, respectively, this translates to 16.2 cm and 26.7 cm of group delay (or phase advance) for every one TEC unit (1 TEC unit = 10^{16} electrons/m²). Daytime and nighttime, as well as solar minimum and maximum ground TEC measurements, vary between 1 and 500 TEC units. Therefore, first-order ionospheric group delay (phase advance) ranges between ~ 0.2 and 80 m for L1 and ~ 0.3 and 130 m for L2.

The first-order ionospheric term, which is about three orders of magnitude larger than higher order terms, can be eliminated by using the "ionospheric free" linear combination, which, based on Eq. (8), is given by

$$\left(\frac{f_1^2}{f_1^2 - f_2^2} \right) P_1 - \left(\frac{f_2^2}{f_1^2 - f_2^2} \right) P_2 = \rho - \frac{s}{f_1 f_2 (f_2 + f_1)} - \frac{r}{f_1^2 f_2^2} \quad (17)$$

As the first-order ionospheric term is eliminated, the dominant ionospheric errors are due to the second- and third-order terms, which are discussed below.

B. Second-Order Effect

The term $B_0 |\cos \theta_B|$ in Eq. (11) represents the absolute value of the component of the \vec{B}_0 field along the line of propagation; therefore, it can be replaced by $|\vec{B}_0 \cdot \vec{k}|$, where (\cdot) represents the inner product and \vec{k} is the unit vector in the direction of propagation.

Consider a station with magnetic colatitude and longitude θ_m and ϕ_m , respectively, observing a satellite with elevation E_m and azimuth A_m , where A_m is measured from magnetic north. Then \vec{k} is given by

$$\vec{k} = -(\cos E_m \sin A_m \hat{X}_m + \cos E_m \cos A_m \hat{Y}_m + \sin E_m \hat{Z}_m) \quad (18)$$

therefore,

$$\left| \mathbf{B}_0 \cdot \vec{k} \right| = B_g \left(\frac{R_E}{r_m} \right)^3 \left| \sin \theta'_m \cos E_m \cos A_m - 2 \cos \theta'_m \sin E_m \right| \quad (19)$$

where θ'_m , r_m are the magnetic colatitude and radial distance of a point along the link, respectively. This term, multiplied by the electron density, is the integrand of Eq. (11), where one must think of r_m and θ'_m as varying along the line of integration. While the exact distribution of electron density along the line of sight is needed to calculate the second-order delay term, a useful approximation can be derived by assuming that the ionosphere consists of a very thin layer at altitude H . Then, the corresponding r_m and θ'_m at the intersection point between the line of sight and the ionospheric layer are given (for $E_m > 10$ deg) by

$$r_m = R_E + H \quad (20a)$$

$$\theta'_m = \theta_m - \frac{H}{R_E \sin E_m} \cos A_m \cos E_m + O\left(\frac{H^2}{R_E^2}\right) \quad (20b)$$

By combining Eqs. (8), (11), and (19), one can approximate the second-order ionospheric group delay (in units of distance) by

$$\begin{aligned} \text{second order ion. group delay} &= 2.61 \times 10^{-18} \lambda^3 \left(\frac{R_E}{r_m} \right)^3 \\ &\times \left| \sin \theta'_m \cos E_m \cos A_m - 2 \cos \theta'_m \sin E_m \right| \text{TEC} \quad (21) \end{aligned}$$

where r_m and θ'_m are given by Eq. (20). Setting H at 300 km and ignoring the factor between the absolute signs, Eq. (21) implies that in the dipole approximation, the second-order ionospheric group delay is on the order of 0.16 mm and 0.33 mm for L1 and L2, respectively, for each TEC unit. The second-order ionospheric phase advance, on the other hand, is one-half of this effect. When forming the ionospheric free linear combination, some cancellation in the second-order term takes place; the residual range error (RRE), which is defined as the difference between the dual-frequency corrected range [left-hand side of Eq. (17)] and the true range, is then on the order of -0.11 mm per TEC unit.

The relations between the magnetic colatitude and longitude, θ_m and ϕ_m , and the geographical colatitude and longitude, θ and ϕ , are given by

$$\begin{aligned} \cos \theta_m &= \sin \delta \cos \beta \sin \theta \cos \phi \\ &+ \sin \delta \sin \beta \sin \theta \sin \phi + \cos \delta \cos \theta \quad (22) \end{aligned}$$

$$\begin{aligned} \tan \phi_m &= \\ &= \frac{-\sin \beta \sin \theta \cos \phi + \cos \beta \sin \theta \sin \phi}{\cos \delta (\cos \beta \sin \theta \cos \phi + \sin \beta \sin \theta \sin \phi) - \sin \delta \cos \theta} \quad (23) \end{aligned}$$

The satellite elevation in local magnetic east-north-vertical coordinates, E_m , is the same as the elevation in local geodetic east-north-vertical coordinates, E . On the other hand, the azimuths in these two coordinates are related through

$$\begin{aligned} A_m &= A + \arccos(\sin \phi \sin \phi_m \cos \delta \cos \beta \\ &+ \cos \phi \cos \phi_m \cos \beta + \sin \phi \cos \phi_m \sin \beta \\ &- \cos \phi \sin \phi_m \cos \delta \sin \beta) \quad (24) \end{aligned}$$

Figure 4 shows the absolute value of the RRE due to the second-order term. This is shown for two stations at different longitudes and latitudes, tracking different GPS satellites, as indicated on the figure. These errors are calculated using the exact integral form of Eq. (11) and assuming the Chapman layer distribution of Fig. 1 and the magnetic field of a tilted dipole, as described above. The angle χ in Eqs. (13) and (14) is determined based on the assumption that the \hat{x} axis (Fig. 3) is pointing toward the sun at 12h UT. The exact calculation, referred to as truth, is compared with an approximation obtained from Eqs. (20)–(24). According to the examples of Fig. 4, the true second-order absolute RRE has an rms value of 1.25 cm, and can be as large as 4 cm at the lowest elevation angle (10 deg). Using the thin-layer model at the 300-km altitude as described above, it is possible to approximate this effect to better than 90 percent on the average. The difference between the truth and the approximation has an average of 0.11 cm and a variance of 0.25 cm. This suggests that a thin-layer model of the ionosphere can be very useful in calibrating the second-order ionospheric effect and therefore improving GPS-user range measurements.

C. Third-Order Effect

Upon examining Eq. (12), one finds that the second term, except during times of very strong magnetic storms, contributes no more than a submillimeter of range error for gigahertz frequencies. Therefore, one must consider the first term, which can be simplified to (in units of meters, kilograms, and seconds)

$$\text{third-order ion. group delay} = 3.0 \times 10^{-31} \lambda^4 \int N^2 dl \quad (25)$$

To get an approximate estimate of the integral of Eq. (25), the authors use the shape parameter η , defined by Brunner and Gu [17] as

$$\eta \equiv \frac{\int N^2 dl}{N_{\max} \int N dl} \quad (26)$$

For a single Chapman layer, η was estimated to be ~ 0.66 and almost independent of elevation [23,17]. Since this ionospheric profile is dominated by a single layer (F2), the authors believe that the shape parameter η in this case will be close to 0.66. Therefore, one can approximate the integral of Eq. (25) by $0.66 \times N_{\max} \times TEC$. For $N_{\max} = 3.0 \times 10^{12} (e/m^3)$ and $TEC = 10^{18} (e/m^2)$ the third-order term is estimated to be ~ 0.86 mm for L1, ~ 2.4 mm for L2, and ~ -0.66 mm for the RRE. A more exact estimate of

the third-order term based on Eq. (12) and the Chapman distribution of Fig. 1 is shown in Fig. 4. In the examples of Fig. 4, the delay ranges between 1 and 4 mm.

VI. Conclusion

The above results are summarized in Table 1, which shows the amount of group delay due to first-, second-, and third-order ionospheric terms in the zenith direction, assuming a zenith $TEC = 10^{18} (e/m^2)$.

In employing a Chapman distribution and a dipole approximation for the magnetic field, it was possible to estimate the higher order ionospheric effects on range and phase measurements. The second-order error can be several centimeters for range as well as phase during daytime, for a year near sunspot maximum. Moreover, since the magnetic field is fixed to the Earth, and the GPS orbit, as seen from a ground station, repeats itself daily (shifted by ~ 4 min per day), the diurnal shape of the second-order error is most likely to repeat its overall structure for several days, at least to the extent that the overall electron density distribution remains unchanged. Such daily repeatable errors in range and phase will be mapped directly into orbital and baseline estimation. This study shows that a rough ionospheric model consisting of a thin shell at 300 km, plus a knowledge of the TEC, allows one to calibrate the second-order term to better than 90 percent. This implies reducing the second-order ionospheric error to less than 2 mm on the average and, therefore, potentially improving orbit determination and baseline solutions.

Acknowledgments

The authors wish to thank Sien Wu and Thomas Yunck of JPL for helpful comments and suggestions on this article.

References

- [1] T. P. Yunck, W. G. Melbourne, and C. L. Thornton, "GPS-based satellite tracking system for precise positioning," *IEEE Tr. Geosci. and Rem. Sensing*, vol. GE-23, pp. 450-457, July 1985.
- [2] S. M. Lichten and J. S. Border, "Strategies for high precision Global Positioning System orbit determination," *J. Geoph. Res.*, vol. 92, no. 10, pp. 12,751-12,762, November 1987.
- [3] T. P. Yunck, S. C. Wu, and J. T. Wu, "Precise Near-Earth Navigation With GPS: A Survey of Techniques," *TDA Progress Report 42-91*, vol. July-September 1987, Jet Propulsion Laboratory, Pasadena, California, pp. 29-45, August 15, 1987.
- [4] S. C. Wu, T. P. Yunck, and C. L. Thornton, "Reduced Dynamic Technique for Precise Orbit Determination of Low Earth Satellites," *J. Guidance, Control and Dynamics*, vol. 14, no. 1, pp. 24-30, January-February 1991.
- [5] S. M. Lichten, "Toward GPS Orbit Accuracy of Tens of Centimeters," *Geoph. Res. Let.*, vol. 17, no. 3, pp. 215-218, March 1990.
- [6] S. M. Lichten, "Estimation and filtering for high-precision GPS positioning application," *Manuscripta Geodaetica*, vol. 15, pp. 159-179, April 1990.
- [7] C. L. Thornton, J. L. Fanselow, and N. A. Renzetti, "GPS-based geodetic measurement systems," *Space Geodesy and Geodynamics*, A. Anderson and A. Cazenave, eds., New York: Academic Press, 1986.
- [8] J. N. Kellogg and T. H. Dixon, "Central and South America GPS geodesy—CASA UNO," *Geoph. Res. Let.*, vol. 17, no. 3, pp. 195-198, March 1990.
- [9] J. T. Freymueller and J. N. Kellogg, "The Extended Tracking Network and Indications of Baseline Precision and Accuracy in the North Andes," *Geoph. Res. Let.*, vol. 17, no. 3, pp. 207-210, March 1990.
- [10] R. P. Malla and S. C. Wu, "GPS Inferred Geocentric Reference Frame for Satellite Positioning and Navigation," *Bul. Geod.*, vol. 63, pp. 263-279, 1989.
- [11] A. P. Freedman, "Measuring Earth Orientation With Global Positioning System," *Bul. Geod.*, vol. 65, pp. 53-65, 1991.
- [12] W. I. Bertiger, J. T. Wu, and S. C. Wu, "Gravity Field Improvement Using GPS Data From TOPEX/POSEIDON: A Covariance Analysis," *J. Geoph. Res.*, vol. 97, no. B2, pp. 1965-1971, February 10, 1992.
- [13] S. M. Lichten, "Precise Estimation of Tropospheric Path Delays With GPS Techniques," *TDA Progress Report 42-100*, vol. October-December 1989, Jet Propulsion Laboratory, Pasadena, California, pp. 1-12, February 15, 1990.
- [14] D. M. Tralli and S. M. Lichten, "Comparison of Kalman filter estimates of zenith atmospheric path delays using the Global Positioning System and very long baseline interferometry," submitted to *Radio Science*.
- [15] D. Coco, "GPS—Satellites of opportunity for ionospheric monitoring," *GPS World*, vol. 2, no. 9, pp. 47-50, October 1991.
- [16] J. J. Spilker, "GPS Signal Structure and Performance Characteristics," *Navigation*, vol. 25, pp. 29-54, 1978.
- [17] F. K. Brunner and M. Gu, "An Improved Model for the Dual Frequency Ionospheric Correction of GPS Observations," *Manuscripta Geodaetica*, vol. 16, no. 3, pp. 205-214, 1991.

- [18] M. Gu and F. K. Brunner, "Theory of the Two Frequency Dispersive Range Correction," *Manuscripta Geodaetica*, vol. 15, pp. 357-361, 1990.
- [19] K. G. Budden, *The Propagation of Radio Waves*, New York: Cambridge Press, 1985.
- [20] A. S. Jursa, ed., *Handbook of Geophysics and the Space Environment*, Air Force Geophysics Laboratory, National Technical Information Services, Springfield, Virginia, 1985.
- [21] C. H. Papas, *Theory of Electromagnetic Wave Propagation*, New York: McGraw-Hill, 1965.
- [22] H. Rishbeth and O. K. Garriott, *Introduction to Ionospheric Physics*, New York: Academic Press, 1969.
- [23] G. K. Hartmann and R. Leitinger, "Range Errors Due to Ionospheric and Tropospheric Effects for Signal Frequencies Above 100 MHz," *Bul. Geod.*, vol. 58, pp. 109-136, 1984.

Table 1. Estimated zenith ionospheric group delay due to $1/f^2$, $1/f^3$, and $1/f^4$ terms, for an arbitrary wavelength λ (microwave region), L1 and L2 frequencies as well as the residual range error with dual-frequency calibration. It is assumed that the zenith $TEC = 10^{18}$ (e/m^2). The phase advance can be read from this table by multiplying each number by = 1, $-1/2$, and $-1/3$ for the $1/f^2$, $1/f^3$, and $1/f^4$ terms, respectively.

Ionospheric expansion term	λ , MKS ^a units	L1	L2	RRE
$1/f^2$	$4.48 \times 10^{-16} \lambda^2 TEC$	16.2 m	26.7 m	0.0
$1/f^3$	$\approx a 2.61 \times 10^{-18} \lambda^3 TEC$ ($0 < a < 2$)	~ 1.6 cm	~ 3.3 cm	~ -1.1 cm
$1/f^4$ ($N_{max} = 3.0 \times 10^{12} e/m^2$)	$\approx 2.0 \times 10^{-31} \lambda^4 N_{max} TEC$	~ 0.86 mm	~ 2.4 mm	~ -0.66 mm
Calibrated $1/f^3$ based on a thin-layer ionospheric model				$\sim 1-2$ mm

^aMeters, kilograms, and seconds.

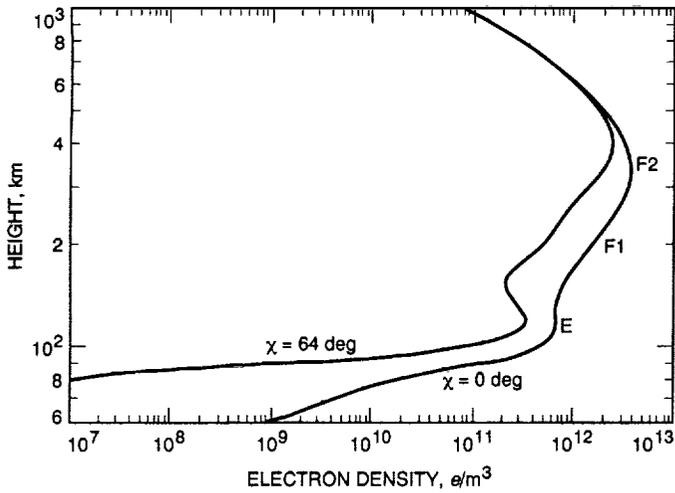


Fig. 1. Ionospheric profile modeled as the sum of three Chapman layers.

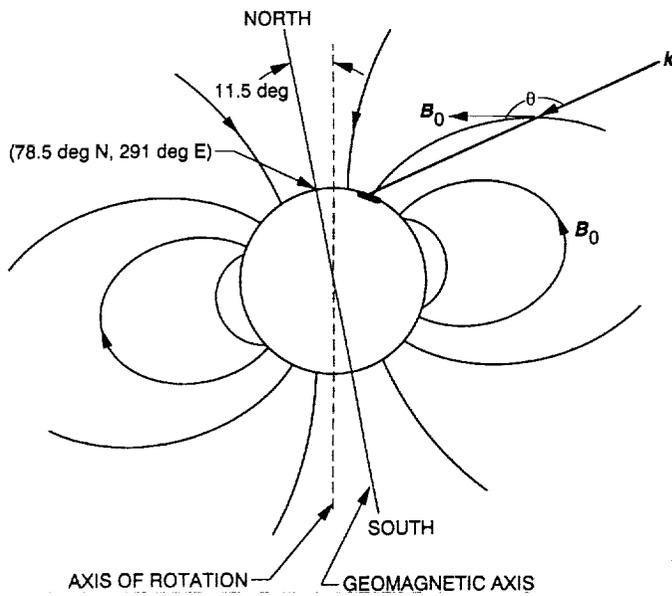


Fig. 2. The Earth's magnetic field modeled as an Earth-centered dipole, aligned along the geomagnetic axis.

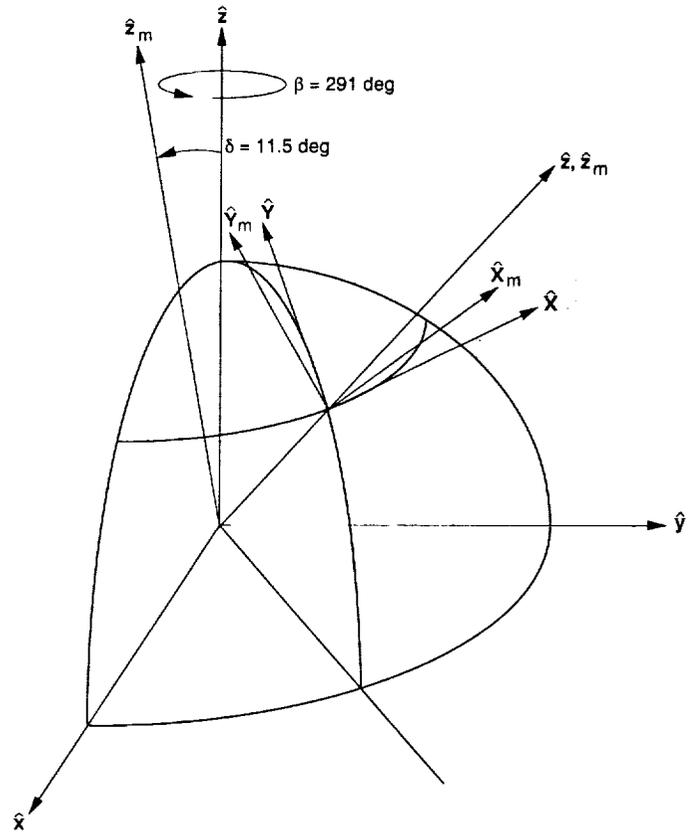


Fig. 3. A graphical illustration of all the frames used in the text. The vectors $\hat{x}, \hat{y},$ and \hat{z} correspond to the geodetic frame; the vectors $\hat{x}_m, \hat{y}_m,$ and \hat{z}_m correspond to the geomagnetic frame; the vectors $\hat{x}_m, \hat{y}_m,$ and \hat{z}_m correspond to geodetic local east, north, and vertical; and the vectors $\hat{x}_m, \hat{y}_m,$ and \hat{z}_m correspond to geomagnetic local east, north, and vertical.

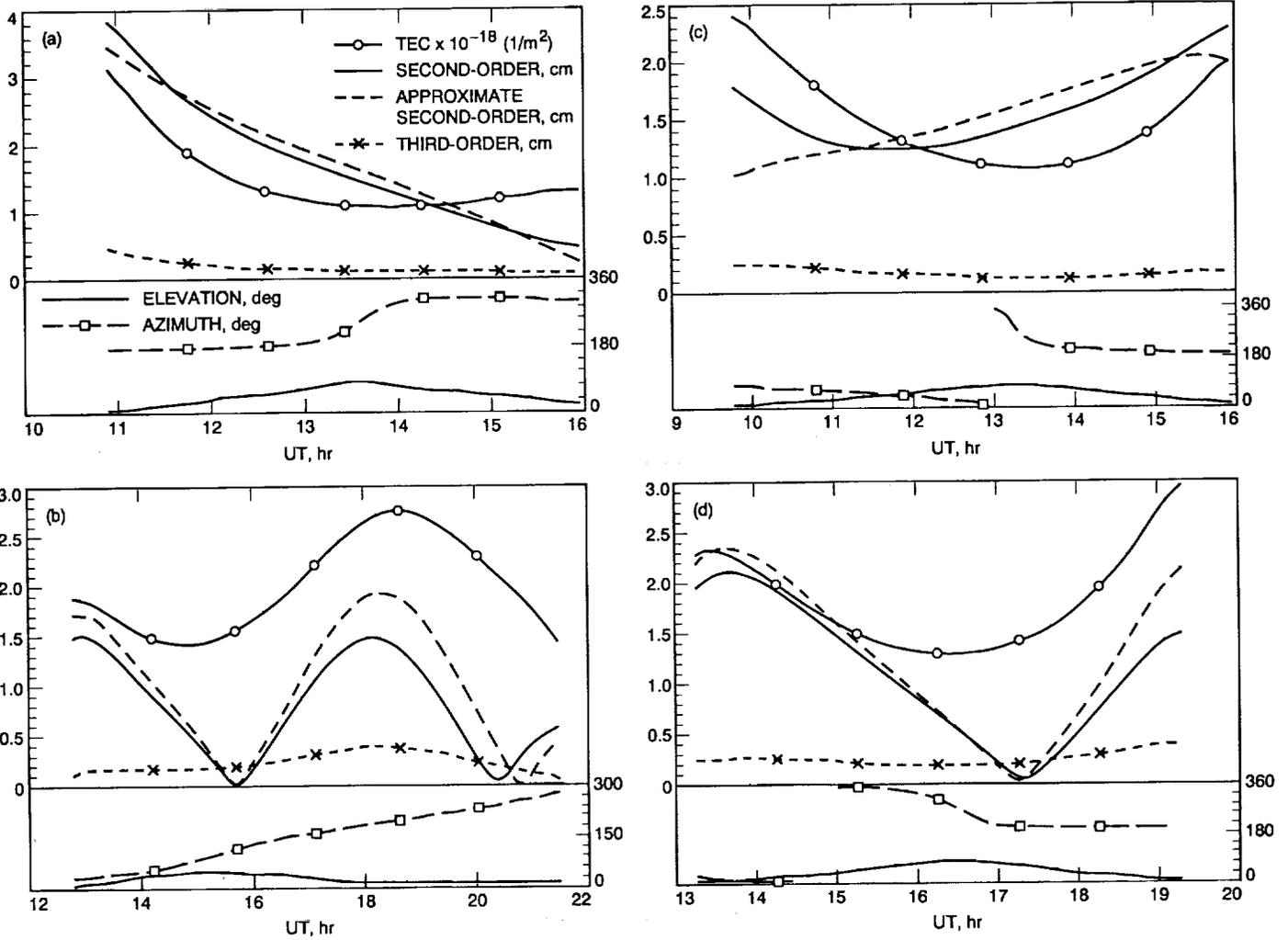


Fig. 4. TEC and absolute second- and third-order ionospheric residual range errors along the line of sight for different GPS-ground receiver pairs. Shown also are the elevation and azimuth of the observed satellite as functions of time: (a) GPS 1, station at 40 deg N latitude, 0 deg longitude; (b) GPS 9, station at 0 deg latitude, 75 deg W longitude; (c) GPS 20, station at 40 deg N latitude, 0 deg longitude; and (d) GPS 16, station at 0 deg latitude, 75 deg W longitude.

58 N 93 - 19421
 120441
 P-14

Evaluation of the Table Mountain Ronchi Telescope for Angular Tracking

G. Lanyi, G. Purcell, and R. Treuhaft
 Tracking Systems and Applications Section

A. Buffington
 Center for Astrophysics and Space Sciences
 University of California at San Diego

The performance of the University of California at San Diego (UCSD) Table Mountain telescope has been evaluated to determine the potential of such an instrument for optical angular tracking. This telescope uses a Ronchi ruling to measure differential positions of stars at the meridian. The Ronchi technique is summarized in this article, and the operational features of the Table Mountain instrument are described. Results from an analytic model, simulations, and actual data are presented that characterize the telescope's current performance. For a star pair of visual magnitude 7, the differential uncertainty of a 5-min observation is about 50 nrad (10 marcsec), and tropospheric fluctuations are the dominant error source. At magnitude 11, the current differential uncertainty is approximately 800 nrad (approximately 170 marcsec). This magnitude is equivalent to that of a 2-W laser with a 0.4-m aperture transmitting to Earth from a spacecraft at Saturn. Photoelectron noise is the dominant error source for stars of visual magnitude 8.5 and fainter. If the photoelectron noise is reduced, ultimately tropospheric fluctuations will be the limiting source of error at an average level of 35 nrad (7 marcsec) for stars approximately 0.25 deg apart. Three near-term strategies are proposed for improving the performance of the telescope to the 10-nrad level: improving the efficiency of the optics, masking background starlight, and averaging tropospheric fluctuations over multiple observations.

I. Introduction

A need for higher data rates and more compact spacecraft hardware has led the Deep Space Network to contemplate using optical communication for the deep-space missions of the next century [1]. As optical communica-

tion is implemented, optical tracking can also be expected to supplement or replace the radio methods now in use. This article explores the applicability of an optical telescope with a Ronchi ruling in the focal plane ("Ronchi telescope") to differential angular tracking of interplanetary spacecraft. In particular, to understand the error

sources that limit both current and ultimate performance, the limiting errors of the UCSD-Table Mountain instrument have been examined using both real and simulated data. The levels of photoelectron and tropospheric fluctuations are analyzed, and analytical models are compared with actual performance over a range of visual magnitudes.

Figure 1 shows schematically the essential features of a Ronchi ruling, which is an optical grating consisting of alternate parallel opaque and transparent lines precisely laid down on a polished glass substrate. Typically there are several line pairs per millimeter, and the widths of the opaque and transparent lines are comparable, if not equal. On a Ronchi telescope the ruling lies in the focal plane, and during a measurement the stars in the field of view move across the ruling at a uniform rate. This motion can be induced either by moving the ruling within the telescope, or by holding the entire telescope fixed and allowing the Earth's rotation to carry the field of view across the ruling. In either case, a detector placed behind the ruling will observe the intensity of starlight modulated periodically by the ruling, as shown on the right side of the figure. If the width of the lines is comparable to the size of the image, the modulation will be continuous and of near-maximum amplitude. A time series of intensity measurements will then contain maximum information about the position of the image on the ruling, in the direction perpendicular to the ruling lines. If two stars are in the field of view simultaneously, then analysis of both time series can give a precise estimate of the difference of the stars' coordinates in that direction. This analysis will be described in more detail in Section II.

There are only three telescopes currently using Ronchi rulings for astrometry. The following list summarizes their distinguishing characteristics:

- (1) Allegheny Observatory [2,3]: Refractor. Telescope tracks to maintain a fixed field of view. Motorized ruling moves across the field in orthogonal directions to determine two coordinates. Masks ("platens") for each field transmit light from only selected stars. Dedicated detector for each star.
- (2) Hipparcos spacecraft, launched August 9, 1989 [4,5]: Reflector. Rotating telescope with slowly variable rotation axis. Optics superimpose two fields 58 deg apart. Fixed ruling rotates with the telescope. Electronic image dissector isolates stars.
- (3) UCSD-Table Mountain [6]: Reflector. Meridian-transit telescope rotates with the Earth, fixed ruling rotates with the telescope. Measures right ascension difference only. Field of view divided into 12 declina-

tion bands with separate detectors. No background masking.

Each design has advantages for particular applications. For spacecraft tracking, a combination of the features listed above would be ideal. As a minimum, such an instrument must be able to determine differential right ascension *and* declination, have a masking capability to block background stars, and be able to reach visual magnitudes, m_v , in the neighborhood of 11. Other desirable—albeit not altogether compatible—features include extreme optical and mechanical stability, a minimum number of parts that move during a measurement, and freedom from optical aberrations over a wide, flat field of view.

Two lines of reasoning support the requirement given above for m_v . First, current projections suggest that spacecraft communication lasers will have apertures of about 0.4 m and transmit 2 W at a wavelength of 0.5 μm .¹ Calculations similar to those described in [7] show that such a laser, transmitting from Saturn, would have an effective m_v of about 11 as seen on Earth. Second, tracking with a Ronchi telescope requires that a reference star be in the same field of view as the spacecraft laser. Suppose that the field of view is 0.5 square degree (as for the Table Mountain telescope) and that observations are required at the point in the ecliptic farthest from the galactic equator, at galactic latitude 60 deg. According to Allen [8], the average density of stars brighter than $m_v = 10$ at that latitude is 4.3 per square degree, and for $m_v = 11$, the density is 11 per square degree. With these parameters, the probability that a random field is empty to $m_v = 10$ is at least 0.12; but at $m_v = 11$, the probability is less than 0.01. Thus, both arguments lead to the conclusion that the telescope must ultimately operate at about $m_v = 11$.

In the material that follows, Section II describes the essential features of the Table Mountain telescope and summarizes the way in which data are collected and analyzed. Section III presents the results of error modeling, simulations, and data analysis that explore the telescope's current and potential capabilities. Finally, Section IV discusses several planned improvements in the design of the telescope that will enable it to approach its ultimate performance.

II. Instrumentation and Data Analysis

The telescope used in these measurements is a Newtonian meridian-circle instrument owned by the University of

¹ J. R. Lesh, personal communication, Communications Systems Research Section, Jet Propulsion Laboratory, May 1990.

California at San Diego and located at JPL's Table Mountain Observatory. Figure 2 is a cross section through the telescope barrel that shows its essential features. Light entering at the right travels down the tube and is reflected from the parabolic primary mirror (M1), which has a diameter of 32 cm and a focal length, F , of 2.443 m. The converging beam then returns to the right along the optical axis to a flat diagonal mirror (M2), which redirects it to the primary focal plane at the Ronchi ruling (R). At this point, the field of view has a diameter of 5 cm, corresponding to 1.2 deg on the sky. Stellar images at the edge of the field are dominated by coma and are about 60 microns long.

On the Ronchi ruling (see Fig. 1), there are 400 line pairs (transparent and opaque) oriented parallel to the declination (north-south) direction in the image. Transparent and opaque lines are equally wide, and the combined width of a line pair is $d = 125$ microns. Thus, each line pair subtends an angle on the sky of

$$\Phi_R = d/F = 10.5545 \text{ seconds of arc} \quad (1)$$

During a measurement, the telescope and ruling remain stationary while the Earth's rotation carries stellar images across the focal plane at a mean angular rate of $\omega = 15.0411$ arcsec/sec at the celestial equator. As an image traverses the ruling perpendicular to the lines, its transmitted light is modulated with a period equal to the time required to cross a line pair. This interval, the Ronchi period, is consequently given by

$$\tau_R = \frac{\Phi_R}{\omega \cos \delta} \quad (2a)$$

where δ is the apparent declination of the star. Substituting the values given above for Φ_R and ω yields

$$\tau_R = 0.70171/\cos \delta \text{ seconds of time} \quad (2b)$$

In essence, the phase of this periodic response of the Ronchi telescope is used to determine the relative right ascension of a star. For example, if two stars at the same declination differ in right ascension by $\Phi_R/2$ seconds of arc, their response functions will be offset by half a Ronchi cycle, and the observed phase difference can be used to deduce the right-ascension difference.

In order to detect the modulated starlight that has passed through the Ronchi ruling, it is convenient to use

a system of transfer optics (labeled M3, L1, L2, and L3 in Fig. 2) to reimagine the star field at the secondary image plane R' . Cylindrical lens L3 produces stellar images that are tightly focused in declination, but diffused along their direction of motion across the field. At R' , a series of 13 razor-edged steel shims extends in the right ascension direction to separate the image into twelve 0.038-deg-wide declination bands or channels. Twelve Plexiglas light pipes convey the light from each channel to a photomultiplier tube. The output current from each tube is then integrated in a capacitor, and the resulting voltage is sampled and digitized at intervals of $\Delta t = 0.075$ sec. Finally, the 12 counts collected for each sampling interval are recorded on a storage device, such as a magnetic disk, for later analysis.

Partitioning the field of view in this way makes it possible to distinguish and analyze separately the instrument's response to as many as 12 stars that are visible simultaneously. In general, of course, there will be several stars in each declination band at any given time, even though all of them may be faint. If the band contains a star to be measured, the cumulative effect of these background stars will influence the telescope's response and may be the dominant source of error. This error consists of two parts: The background stars increase the level of stochastic photoelectron noise, and also introduce systematic offsets in the estimated position as their response functions interfere with the response function of the star being measured. The background problem is discussed further in Sections III and IV.

As the foregoing discussion makes clear, the Ronchi telescope's response to a single star is a time series of photomultiplier counts (shown in Fig. 1) that rises and falls as the stellar image passes in turn across transparent and opaque lines of the Ronchi ruling. Because the width of the lines is comparable to the maximum coma-broadened size of a stellar image, the amplitude of the modulation is nearly 100 percent. The average number of sample points in a Ronchi period, N_R , is simply

$$N_R = \tau_R/\Delta t = 9.3562/\cos \delta \quad (3)$$

Because N_R provides the most precise determination of focal length, the experimentally determined constant 9.3562 in Eq. (3) is used to calculate τ_R and Φ_R .

Optical stability, and mechanical and optical simplicity, were the primary considerations in the design of the Table Mountain telescope. Nothing moves during an observing session, and the absence of lenses in the primary optics minimizes chromatic aberration. The images do,

however, have the coma characteristic of a single parabolic mirror. As a result, the response function is not strictly periodic, and systematic errors can arise in the comparison of Ronchi phases measured at different points in the field of view.

Coma in the Table Mountain instrument is accounted for during data analysis. Buffington and Geller [6] argued that even in the presence of coma, the centroid of the image does move uniformly across the ruling. Furthermore, the centroid of a single peak in the response time series occurs when the centroid of the image crosses the midline of one of the transparent strips on the ruling. Hence, the problem of determining the phase of the Ronchi response function reduces to that of determining the centroid times of the response time series. That is, the analysis software computes for each peak in the response the centroid time,

$$\langle t \rangle_k = \frac{\sum_{i=1}^n t_i A_i}{\sum_{i=1}^n A_i} \quad (4)$$

where $\langle t \rangle_k$ is the time of the k th centroid, t_i is the i th sample time in the interval, and the values of A_i are the corresponding modulated intensities after removal of a background level that varies with time. The sum is taken over n data points in a particular peak.

In the absence of noise, the centroids of successive Ronchi cycles recur at equal intervals of τ_R , so that the Ronchi phase of the k th centroid can be defined as

$$\phi_k = (\langle t \rangle_k - t_{ref}) / \tau_R - k \quad (5a)$$

where ϕ_k is in cycles and t_{ref} is a reference time common to all stars. Ronchi phase can be expressed as an angular offset in the right ascension direction simply by multiplying ϕ_k by the conversion factor Φ_R

$$\alpha_k = \Phi_R \phi_k \quad (5b)$$

Ideally, ϕ_k is constant for a particular star, but in the presence of perturbations (photoelectron noise, background stars, tropospheric refraction, and so on) it fluctuates. However, because the tropospheric fluctuations are correlated for stars separated by small angles, part of this error source cancels in a differential measurement between stars in different declination bands. Figure 3 shows the time series of Ronchi phase for a pair of magnitude-7 stars separated by about 23 minutes of arc. In the figure, the solid and dashed lines represent the phases of the

two stars, and dots show the difference. To show the correlation of the tropospheric fluctuations more clearly, the mean phases of the two stars have been made equal. Because the troposphere is not the only source of error, and because the troposphere itself is not perfectly correlated, the correlation coefficient of the two time series is only about +0.5. In the following section, data on observed tropospheric fluctuations are presented, and the variation of photoelectron noise as a function of stellar brightness is discussed in detail.

III. Results

The goal in assessing the Table Mountain Ronchi telescope has been to evaluate the errors limiting its current performance and to determine how much that performance can be improved within the constraints imposed by its basic design. Ultimately, it will be necessary to decide whether such an instrument can track interplanetary spacecraft with the required accuracy. As discussed above, the two limiting error sources are photoelectron noise and tropospheric fluctuations.

In order to evaluate these two sources of error quantitatively, a combination of theoretical analysis, simulated data, and actual data has been used. Figures 4 and 5 summarize these results for differential observations of stars ranging from $m_v = 4.5$ to 12. In these figures, observed and calculated errors, in nanoradians, are plotted as a function of m_v .

Along the diagonal in Fig. 4, the solid line represents an analytic model of the component of angular error induced by photoelectron noise for the current instrument. The Appendix gives the derivation of this error for a single centroid measurement. To extrapolate the single-centroid calculation to a full-length (400-centroid) differential observation, a reasonable assumption was made that photoelectron noise on different centroids, or in different detectors, is statistically independent. Thus, the error on an average of 400 centroids is reduced by a factor of $\sqrt{400} = 20$ relative to a single centroid; and the error on a differential measurement is $\sqrt{2}$ times the error on a single-star measurement. Under this assumption, the plotted curve represents Eq. (A-9) divided by $\sqrt{200}$; that is,

$$\sigma_\alpha = (2089/A)(1 + 1.9604 \times 10^{-3}A)^{1/2} \quad (6)$$

where σ_α , the angular uncertainty, is expressed in nanoradians, $A = 10^{0.4(12-m_v)}$ is proportional to the star's brightness [see Eq. (A-8)], and m_v is its visual magnitude.

Note that the calculation of Eq. (6) (and the plotted line) assumes that the star is close to the celestial equator, where there are about nine sample points in a Ronchi period. Away from the equator, σ_α varies as $\sqrt{\cos \delta}$, where δ is the declination [see Eqs. (3) and (A-7)].

The six diamonds plotted in Fig. 4 represent actual differential measurements on six pairs of stars. Table 1 lists the stars in each pair, in order from left to right on the plot, along with their visual magnitudes and angular separations. For each pair, the plotted magnitude is an effective magnitude that accounts for the difference in magnitude of the two stars and the actual spectral response of the photomultiplier tubes.

The results for pairs 1 through 4 are derived from data collected on seven nights between May 23 and June 2, 1990. For each pair, the rms variation over the seven nights of the single-centroid differenced Ronchi phases was first calculated. Then, the single-centroid standard deviations were scaled to the length of a full observation by dividing by 20, as described above. This procedure implicitly assumes that the tropospheric, as well as the photoelectron, fluctuations scale with time as $t^{-1/2}$. Lindegren's semi-theoretical estimate of differential tropospheric fluctuations [9] implies that this relation is correct, although his results do not strictly apply to the combination of angular separation (3 to 25 arcmin) and time interval (τ_R approximately 0.7 sec) applicable here. However, for undifferenced measurements, Lindegren [9] and his references expect the fluctuations to scale like $t^{-\beta}$, where β is between 1/6 and 2/5.

At the right in Fig. 4, the points for the two brightest pairs were obtained indirectly, by extrapolating the data given by Buffington and Geller [6] in their Fig. 5. That figure, based on five nights of measurements made from June 9 to June 13, 1989, shows angular precision as a function of integration time. Here their results have been extended to an integration time of 300 sec, assuming the $t^{-1/2}$ dependence suggested by the plot for shorter integration times.

Examination of the points shows that the weakest pairs lie near the theoretical photoelectron noise curve, but that, for brighter stars, the predicted photoelectron noise increasingly underestimates the actual angular uncertainty. Thus, it appears that for stars weaker than visual magnitude 8.5, photoelectron noise is the dominant error source. At magnitude 7.5, photoelectron and tropospheric fluctuations contribute about equally. For stars brighter than $m_v = 7$, the troposphere dominates the error budget. In Fig. 4, the rightmost points suggest that the error has

nearly reached the asymptotic level of the tropospheric contribution alone. The horizontal line at about 35 nrad indicates an empirical estimate of that limit (for a single measurement) based on the data shown in the figure. Of course, this limit is merely representative. It depends on angular separation and on the time and place of the measurements.

Finally, the sloping line at the lower left in Fig. 4 shows the reduction in photoelectron noise expected to result from two improvements in telescope design discussed more fully in Section IV. First, masks are used to block background stars and remove their contribution to the noise; and second, the efficiency of the transfer optics and light pipes is improved by a factor of 12, so that the signal-to-noise ratio (SNR) increases by a factor of $\sqrt{12}$. With these improvements, the troposphere will dominate the error budget even at $m_v = 12$.

Figure 5 shows again the computed angular error due to photoelectron noise, and compares it with the results of tests in which the existing analysis software was used to process simulated data. In these tests, 1200 undifferenced 400-centroid observations were simulated at each of 16 visual magnitudes ranging from 4.5 to 12 in steps of 0.5. Gaussian noise representing photoelectron (but not tropospheric) fluctuations was added to a sinusoidal approximation of the Ronchi response function, and the results were written to a file in the same format as real data. The analysis software was then used in the usual way to compute an angular coordinate for each observation. Finally, each observation was assigned an error equal to the difference between the computed coordinate and the model coordinate used to generate the sinusoid. The plotted values along the solid line show the standard deviation of the 1200 errors, multiplied by $\sqrt{2}$ to account for differencing.

Except for the brightest and faintest stars, the predicted and simulated uncertainties agree remarkably well. It is not yet well understood why the simulations perform better than the model for stars of visual magnitude 12. The breakdown of the assumptions underlying Eq. (A-6) certainly plays a role, however. At the bright end of Fig. 5, it is suspected that an undiagnosed algorithm error is limiting the uncertainty derived from the simulated data at the 6-nrad level. If so, the simulation curve will agree better with the analytic model when the error is corrected.

Figure 6 applies the results shown in Fig. 4 to the special case of spacecraft tracking. As stated in Section I, a spacecraft laser with nominal characteristics at the distance of Saturn would generate a response in a Ronchi telescope comparable to that of a magnitude-11 star. Assuming a reference star of the same magnitude, the figure

shows the expected angular accuracy for several situations. At the left is shown the performance of the Table Mountain instrument, both in its current configuration and with the improvements mentioned above. On the right are the tropospheric errors for a single 5-min observation (taken from Fig. 4) and for the average of 25 statistically independent observations. With the expected improvements and multiple measurements, the estimated accuracy is adequate for near-term research and development demonstrations of optical astrometry.

IV. Conclusion

Very long baseline radio interferometric astrometry can now achieve an angular accuracy of 1 nrad or better [10]. Optical tracking methods must therefore strive toward a comparable goal. For ground-based systems in the near term, 10 nrad is a reasonable target. From the discussion in Section III, it follows that the Table Mountain telescope requires the improvements summarized below before it can deliver the desired performance.

For bright stars one or two tenths of a degree apart, tropospheric fluctuations typically limit differential accuracy to 30 or 40 nrad for a single 5-min measurement. Tropospheric error cannot be controlled, but it can be managed to some extent by observing at high altitude and using star pairs separated by small angles. It can also be reduced by averaging together several measurements, as indicated in Fig. 6.

For stars fainter than $m_v = 8.5$, photoelectron noise dominates the error budget. As pointed out in Section I, calculations of the apparent brightness of spacecraft lasers and of the number of observable stars lead to the conclusion that differential tracking measurements will have to rely on stars with m_v approximately 11. For such a pair (see Fig. 6), the differential angular uncertainty of a 5-min measurement with the current system is about 830 nrad.

Steps are already being taken that will reduce the photoelectron noise on stars of visual magnitude 11 to a level comparable to the tropospheric noise. An all-mirror Offner system [11] will replace the current lens system of transfer optics, increasing the telescope's optical transmission by a factor of 4. A further change in the way the light pipes are connected to the photomultipliers is expected to increase the number of photons that reach the detectors by another factor of 3. Since photoelectron noise varies as the square root of the incident intensity, this 12-fold increase in optical efficiency will increase the SNR by a factor of $\sqrt{12}$ and decrease the uncertainty of a measurement by the same amount.

Another factor that increases photoelectron noise is background light from fainter stars in each declination band. Typically, the total background in a band is comparable to the light from a magnitude-7 star. As mentioned in Section II, this background introduces both random and systematic errors into the estimated coordinates. To remove both kinds of errors, a masking device that will either block or ignore background light is being added. As a preliminary implementation, mechanical masks are being designed on at least two of the channels. Each mask will be an opaque strip that covers one declination band and contains a pinhole to allow only the light from a single star to pass. As the star crosses the field of view, a computer-controlled drive mechanism will move the mask to keep the star centered on the pinhole.

A much more versatile electronic masking system would use a sensitive charge-coupled device (CCD) to replace both the mechanical masks and the photomultiplier tubes. Such a system would retain the Ronchi ruling to modulate the starlight and would use the CCD as a masked detector. Only those CCD pixels containing the desired image would be processed, while those containing background would be discarded. A more radical departure from the current design would use the CCD not only to replace the masks and photomultipliers, but also as a metric device to replace the ruling itself. This pure CCD design may be subject to systematic errors that are difficult to control, however [12,13]. Although CCD's have been used in astrometry for over 10 years [14,15], the design envisioned here presents new challenges. In particular, it requires CCD's that are larger and can be read out faster than those now readily available. Thus, the CCD concept would be developed gradually only after a successful demonstration of mechanical masking.

Figure 4 shows the reduction in photoelectron error that is expected after the implementation of both background masking and improved transfer optics. With these improvements, the photoelectron error will be reduced below the tropospheric limit even for stars of visual magnitude 12.

Finally, tracking applications will require the measurement of spacecraft declination as well as right ascension. This capability can be added to the existing instrument in several ways. For example, stars would move obliquely across the ruling during observations made before or after transit. A pair or series of such observations could be combined to give a differential measurement of both right ascension and declination. Of course, this option would require modification of the telescope to allow nontransit measurements. Another approach would use a chevron

ruling with separate sections of lines oriented at ± 45 deg with respect to the vertical. Still another possibility is to move the ruling itself in the declination direction, as Gatewood's instrument does [2,3]. Some of these methods are affected by differential atmospheric refraction, however, and it is still unclear what will be the best approach for two-dimensional measurements.

In summary, the following modifications would prepare the Table Mountain Ronchi telescope for a demonstration of its ability to track objects as faint as $m_v = 11$:

- (1) Install the transfer optics now being developed, so as to improve the photoelectron SNR.
- (2) Reconfigure the interface between the light pipes and the photomultipliers, which would also improve the photoelectron SNR.
- (3) Design and install computer-controlled masks for at least two declination bands to eliminate the photo-

electron noise and systematic errors caused by background stars.

The following two items offer some potential for improvement, but their feasibility has not yet been studied:

- (1) Replace the photomultiplier tube assembly with a CCD (positioned so the image is slightly out of focus) to investigate the use of CCD's for both detection and masking.
- (2) Install a Ronchi ruling with a chevron pattern for simultaneous measurement of right ascension and declination.

These improvements can be implemented on the current instrument with a modest investment, and they will make it possible to assess the applicability of the Ronchi technique to optical tracking.

References

- [1] J. R. Lesh, L. J. Deutsch, and W. J. Weber, "A Plan for the Development and Demonstration of Optical Communications for Deep Space," *TDA Progress Report 42-103*, vol. July-September 1990, Jet Propulsion Laboratory, Pasadena, California, pp. 97-109, November 15, 1990.
- [2] G. Gatewood, L. Breakiron, R. Goebel, S. Kipp, J. Russell, and J. Stein, "On the Astrometric Detection of Neighboring Planetary Systems," *Icarus*, vol. 41, pp. 205-231, February 1980.
- [3] G. D. Gatewood, "The Multichannel Astrometric Photometer and Atmospheric Limitations in the Measurement of Relative Positions," *Astron. Journal*, vol. 94, pp. 213-224, July 1987.
- [4] J. Kovalevsky, "Prospects for Space Stellar Astrometry," *Space Science Review*, vol. 39, nos. 1/2, pp. 1-63, September/October 1984.
- [5] L. Lindegren, "Hipparcos Data Reduction Overview," *Adv. Space Research*, vol. 11, no. 2, pp. 25-34, 1991.
- [6] A. Buffington and M. E. Geller, "A Photoelectric Astrometric Telescope Using a Ronchi Ruling," *Publications of the Astronomical Society of the Pacific*, vol. 102, no. 648, pp. 200-211, February 1990.
- [7] B. L. Schumaker, "Apparent Brightness of Stars and Lasers," *TDA Progress Report 42-93*, vol. January-March 1988, Jet Propulsion Laboratory, Pasadena, California, pp. 111-130, May 15, 1988.
- [8] C. W. Allen, *Astrophysical Quantities*, Third Edition, London: The Athlone Press, 1973.

- [9] L. Lindegren, "Atmospheric Limitations of Narrow-Field Optical Astrometry," *Astron. and Astrophysics*, vol. 89, nos. 1/2, pp. 41-47, September 1980.
- [10] R. N. Treuhaft and S. T. Lowe, "A Measurement of Planetary Relativistic Deflection," *Astron. Journal*, vol. 102, pp. 1879-1888, November 1991.
- [11] A. Offner, "New Concepts in Projection Mask Aligners," *Optical Engineering*, vol. 14, no. 2, pp. 130-132, March/April 1975.
- [12] A. Buffington, H. S. Hudson, and C. H. Booth, "A Laboratory Measurement of CCD Photometric and Dimensional Stability," *Publications of the Astronomical Society of the Pacific*, vol. 102, pp. 688-697, June 1990.
- [13] A. Buffington, C. H. Booth, and H. S. Hudson, "Using Image Area to Control CCD Systematic Errors in Spaceborne Photometric and Astrometric Time Series Measurements," *Publications of the Astronomical Society of the Pacific*, vol. 103, pp. 685-693, July 1991.
- [14] D. G. Monet and C. C. Dahn, "CCD Astrometry. I. Preliminary Results from the KPNO 4-m/CCD Parallax Program," *Astron. Journal*, vol. 88, pp. 1489-1507, October 1983.
- [15] D. G. Monet, "Recent Advances in Optical Astrometry," *Ann. Rev. Astron. Astrophysics*, vol. 26, pp. 413-440, 1988.

Table 1. Star pairs shown in Fig. 4.

Pair	Stars (SAO identification)	Visual magnitudes	Separation, arcmin
1	122716, 122746	7.46, 8.20	24
2	122735, 122738	6.93, 8.40	13
3	122723, 122709	6.66, 7.60	14
4	122723, 122715	6.66, 7.22	23
5	70287, 70289	6.32, 6.49	3
6	101145, 101137	3.86, 5.91	13

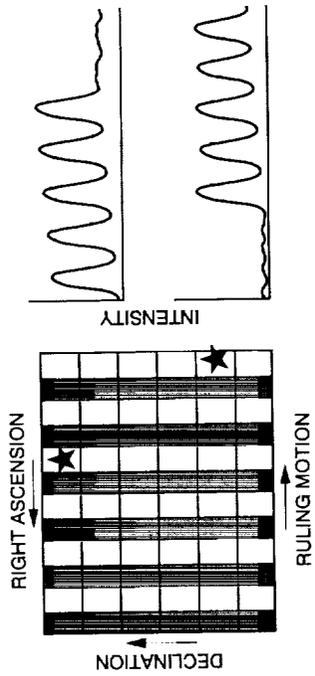


Fig. 1. Ronchi ruling: principle of operation.

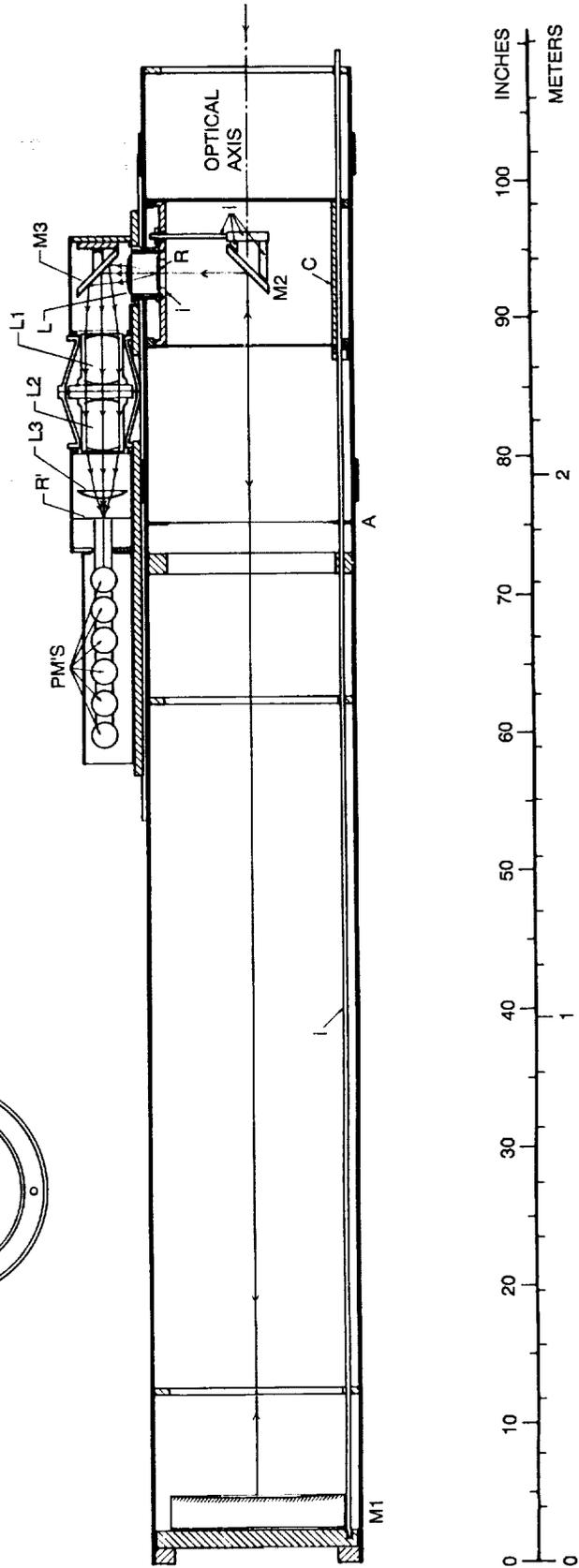
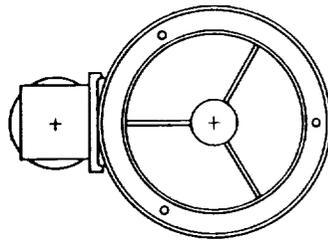


Fig. 2. Telescope configuration.

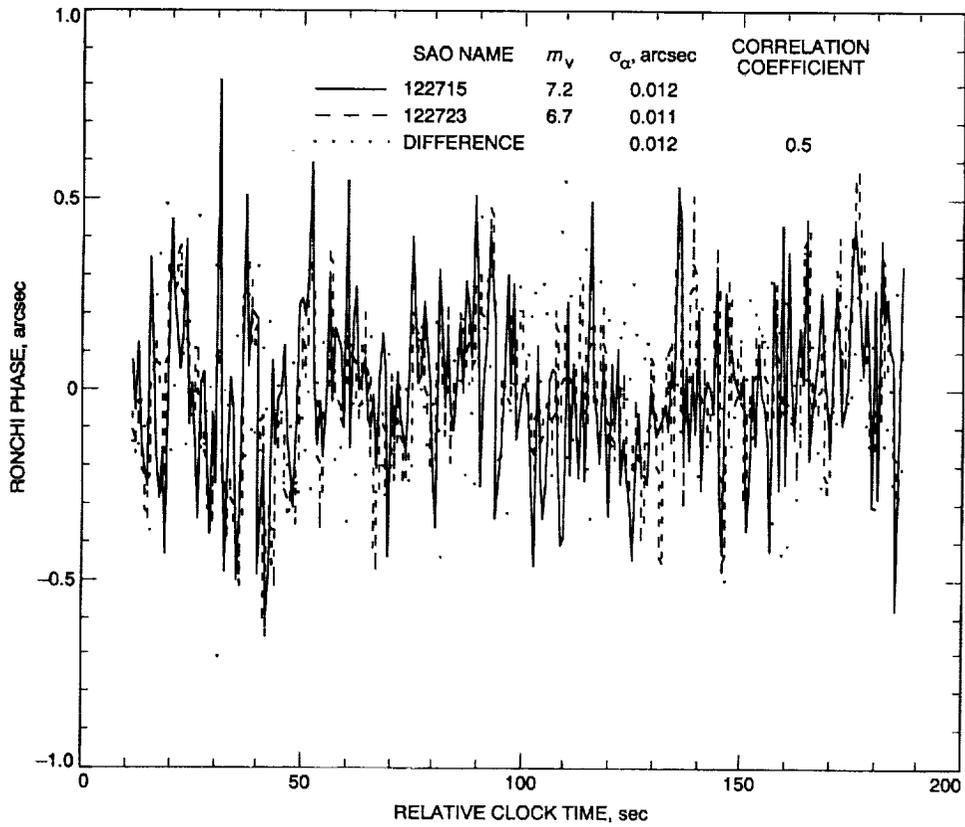


Fig. 3. Observed Ronchi phase.

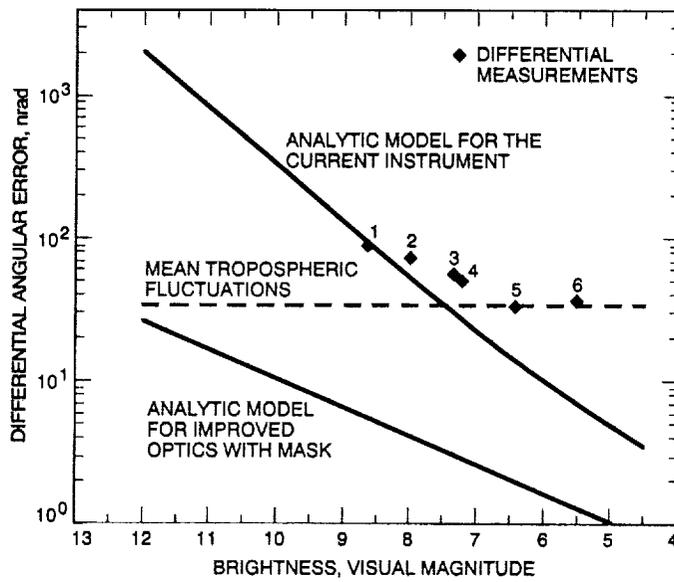


Fig. 4. Differential angular performance.

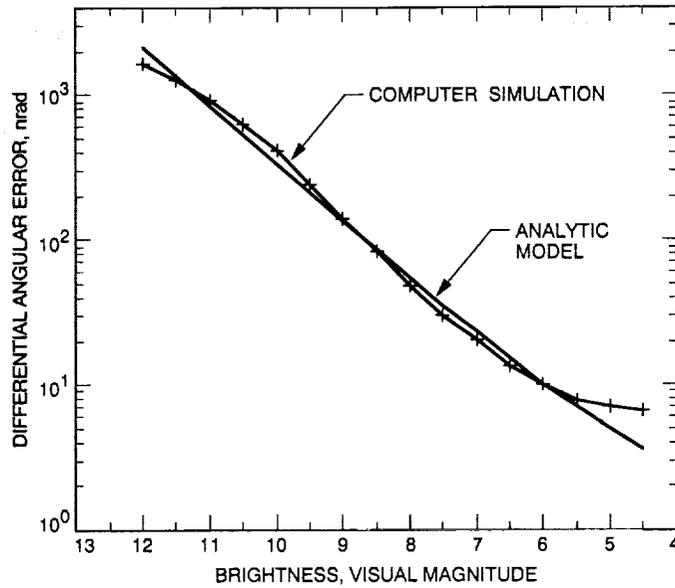


Fig. 5. Comparison of analytic error model with computer simulation. A relative uncertainty of 0.02 is associated with each simulated angular error value.

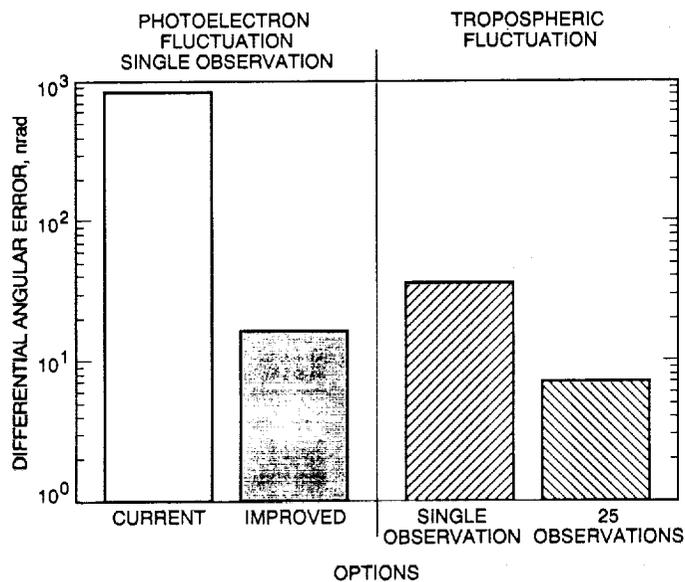


Fig. 6. Error budget for visual magnitude 11.

Appendix

Error in Estimate of Ronchi Centroid Caused by Photoelectron Fluctuations

In the Table Mountain Ronchi telescope, 12 analog-to-digital converters (ADC's) process the signals from 12 corresponding photomultipliers that have unequal gains. Where numerical constants are given below, they refer to channel 3. For other channels, the constants must be scaled by the appropriate gain ratio.

Since the response of each photomultiplier, A , is proportional to the number of photoelectrons, N , and N is proportional to the number of incident photons, the output from each ADC can be written as

$$A = KN \quad (\text{A-1})$$

and is proportional to the brightness of the star being observed.

In the absence of a bright star in the field, there remains an average background level, A_B , due to faint stars. Even if the field of view remains fixed, this background has a stochastic fluctuation level (standard deviation), σ_B . From measurements of "empty" fields, the numerical values of these constants are found to be $A_B \approx 100$ and $\sigma_B \approx 3$. Since the standard deviation of the number of photoelectrons during any interval is the square root of the mean number during that interval (Poisson statistics), the fluctuation of Eq. (A-1) can be written as

$$\sigma = K\sqrt{N} \quad (\text{A-2})$$

Thus, the value of K can be determined from the measured values of A_B and σ_B

$$K = \sigma_B^2/A_B \quad (\text{A-3})$$

and from the values above, $K \approx 0.09$.

When the number of sample points in a Ronchi cycle is odd, $n = 2l + 1$, the location of the centroid of the response function can be written as

$$x_c = (\Phi_R/n) \frac{\sum_{i=-l}^l i A_i}{\sum_{i=-l}^l A_i} \quad (\text{A-4})$$

where the index i ranges over the points in a single Ronchi cycle, and x_c is expressed in the same units as Φ_R , the angle subtended on the sky by one Ronchi line pair [see Eq. (1)]. Note that the values of A_i in Eq. (A-4) are the amplitudes of the response time series of the single star being measured, after subtraction of the background level.

If the fluctuation of the amplitude, σ_i , is much smaller than the amplitude, A_i , then the variance of x_c is approximately

$$\sigma_{x_c}^2 = \frac{(\Phi_R/n)^2 \sum_{i=-l}^l i^2 \sigma_i^2}{\left(\sum_{i=-l}^l A_i\right)^2} \quad (\text{A-5})$$

where $\sigma_i = K\sqrt{N_B + N_{S_i}}$, and the subscripts B and S refer to the background and source, respectively. The value of N_B is more or less constant, but N_{S_i} varies directly with A_i . If the size of the stellar image on the Ronchi ruling is comparable to, or larger than, the line spacing, then the response function is roughly sinusoidal, and the response time series can be approximated as

$$A_i = \frac{A_S [1 + \cos(2\pi i/n)]}{2} \quad (\text{A-6})$$

Using Eqs. (A-2), (A-3), and (A-6) in Eq. (A-5) and applying further numerical approximations, one obtains

$$\sigma_{x_c} = \frac{\Phi_R}{\sqrt{3}n} \sqrt{(\sigma_B/A_S)^2 + \frac{1}{2}(1 - 6/\pi^2)(\sigma_B/A_B)(\sigma_B/A_S)} \quad (\text{A-7})$$

Measurements show that for a star of $m_v = 7$, $A_S \approx 100$, so that the relationship between signal amplitude and stellar magnitude is approximately

$$A = 10^{0.4(12 - m_v)} \quad (\text{A-8})$$

The approximations used in deriving Eq. (A-7) become significant for small values of n and A_S . For $n \geq 9$ (the smallest possible value) and $A_S \geq 30$ (corresponding to $m_v \approx 8$), the error is at most a few percent.

When one substitutes $n = 9$ for stars near the celestial equator, the numerical values given above for A_B and

σ_B , and the value of Φ_R from Eq. (1), then Eq. (A-7) becomes

$$\sigma_{x_c} = 2.03\sqrt{(3/A_S)^2 + (0.2)0.03(3/A_S)} \quad (\text{A-9})$$

where σ_{x_c} , the uncertainty in the location of the star caused by the photoelectron noise, is expressed in seconds of arc.

If the background is removed, only the second term in Eq. (A-9) remains

$$\sigma_{x_c} = 2.03\sqrt{(0.2)0.03(3/A_S)} \quad (\text{A-10})$$

Eq. (A-10) is correct within a few percent for stars brighter than $m_v \approx 11$. For stars fainter than the specified limits, Eqs. (A-7) and (A-10) both overestimate the actual uncertainties.

59 N93-19422

128442

P-10

Deep-Space Navigation Applications of Improved Ground-Based Optical Astrometry

G. W. Null, W. M. Owen, Jr., and S. P. Synnott
Navigation Systems Section

Improvements in ground-based optical astrometry will eventually be required for navigation of interplanetary spacecraft when these spacecraft communicate at optical wavelengths. Although such spacecraft may be some years off, preliminary versions of the astrometric technology can also be used to obtain navigational improvements for the Galileo and Cassini missions. This article describes a technology-development and observational program to accomplish this, including a cooperative effort with U.S. Naval Observatory Flagstaff Station. For Galileo, Earth-based astrometry of Jupiter's Galilean satellites may improve their ephemeris accuracy by a factor of 3 to 6. This would reduce the requirement for onboard optical navigation pictures, so that more of the data transmission capability (currently limited by high-gain antenna deployment problems) can be used for science data. Also, observations of European Space Agency (ESA) Hipparcos stars with asteroid 243 Ida may provide significantly improved navigation accuracy for a planned August 1993 Galileo spacecraft encounter.

I. Introduction

There is an active technology development effort [1,2] at JPL, supporting possible implementation of a ground-based Deep Space Optical Reception Antenna (DSORA), consisting of a 10-m segmented receiving mirror for the downlink and a smaller, roughly 1-m, uplink telescope. This system would provide laser communications between the DSN and interplanetary spacecraft. Although the primary purpose of the DSORA would be to improve deep-space downlink communication, optical tracking systems could also provide new capabilities for interplanetary navigation.

For example, it may eventually be possible to directly image the spacecraft relative to solar system target bodies,

thus providing ground-based optical navigation roughly comparable to the current onboard optical navigation capability, and with approximately the same linear accuracy at the distance of Jupiter. This will eventually require significant improvements in existing astrometric instruments and techniques. Since laser-emitting spacecraft will probably not be available in the 1990s, it is necessary to develop astrometric systems for observing these spacecraft without initially being able to observe them. Fortunately, a suitable astrometric replacement is to observe natural bodies (satellites or asteroids) together with background stars.

As the instrument development proceeds, the improved instruments can also provide improved mission target location accuracy for conventional radio metric missions, such

as Galileo. Such improvements are especially important since, as will be discussed, target location is the limiting error source for the most critical portions of these missions. Two potentially significant improvements are identified in this article: supporting navigation of the Galileo-Jupiter orbit tour and a possible Galileo flyby of asteroid Ida. Similar benefits may be possible in the future during the Cassini-Saturn orbit tour.

Most long-term development options are expensive and technically difficult, but some simpler near-term options have been identified for filled-aperture instruments with narrow (<1 deg) fields. These near-term options and their potential navigation benefits are the subject of the present article. It will focus on optical astrometry for target location, i.e., measurement of the angular sky-plane position of solar system objects relative to one another or to background stars.

One option, which is being actively pursued in cooperation with the United States Naval Observatory (USNO) Flagstaff Station (NOFS), is to observe Jupiter's Galilean satellites with NOFS's 1.55-m astrometric reflector and a modern, large format 2048×2048 charge-coupled device (CCD) detector. Although the 11-arcmin field is unusually large for a CCD detector, it is quite small as compared with a conventional photographic field.

The near-term goal is to achieve per-night accuracy of roughly 50 nrad (≈ 0.01 arcsec), although useful results can still be achieved with larger errors. This goal represents an appropriate trade-off between mission needs and near-term observational limitations, as will be explained later. If this observational goal can be achieved, then the Galileo-Jupiter tour navigation performance can be significantly improved.

The key technology challenge with an 11-arcmin CCD field¹ is to accurately determine the instrument scale (for CCD's, in arcsec/pixel measured on the focal plane), even though, within this field, existing star catalogs do not provide an adequate number of stars with accurately known positions. Observations of star fields are currently being acquired and analyzed by NOFS so that candidate scale-determination techniques can be assessed. These techniques are discussed in Section IV.A.

Five major sections are included in this article: Introduction, Navigation Overview, Instrument Overview, Galilean Satellite CCD Observation Techniques, and Summary and Conclusions.

¹ D. G. Monet, personal communication, NOFS, February 12, 1992.

II. Navigation Overview

A. Galileo-Jupiter Orbit Tour

The Galileo spacecraft will be injected into orbit about Jupiter in late 1995, followed by a series of close encounters with Jupiter's Galilean satellites. Unfortunately, the spacecraft high-gain antenna did not fully deploy and currently is completely unusable. The low-gain antenna is available, but the data rate from Jupiter allows transmission of only a few full-field CCD science or onboard optical navigation (OPNAV) pictures per encounter (there is one encounter per 14-28 day orbit). Obviously, it would be beneficial to take more science and fewer OPNAV images, provided that mission navigation requirements can still be met.

Before the antenna deployment problem, roughly 30 onboard OPNAV pictures were planned for each satellite flyby, but now, even with data compression, there is a strong benefit if the number of pictures can be significantly reduced. Also, having both OPNAV and ground-based optical information can provide increased navigation reliability, and can also provide a quick, accurate three-dimensional target-location position fix by combining angular observations from two different lines of sight.

Close-up OPNAV satellite images can provide an accuracy of about 15 km (about 25-nrad, geocentric). If the 50-100 km a priori position accuracy of the Galilean satellites can be improved to the OPNAV accuracy levels, then the originally planned spacecraft navigation accuracy can be achieved with just a few OPNAV pictures to locate the spacecraft relative to the already well-known position of the target satellite.

As will be discussed later, the orbit improvements would be made with ground-based intersatellite observations. Figure 1 shows a schematic representation of typical observing geometry for these observations. Figure 2 shows an actual CCD frame containing four exposures of the Galilean satellites, taken with the NOFS 1.55-m telescope. The shutter was closed and the pointing was offset between each exposure. The leftmost satellite of the leftmost exposure is outside the field; otherwise, Jupiter and all four satellites are imaged for each exposure, and they are nearly collinear, roughly in the ecliptic plane. The satellites appear at a shallow, roughly 30-deg angle from the horizontal; the leftmost three satellites are somewhat below the Jupiter location.

The outermost Galilean satellite (Callisto) has a longer period and larger intersatellite angular separations than

the other Galilean satellites and, therefore, its data coverage, particularly for eclipses [3] and mutual events [4], is significantly less complete. Thus, it is not surprising that Callisto actually has the least accurate orbit of any Galilean satellite, with a longitude standard error² of about 90 km, and there is a high priority on improving the Callisto orbit. Since the angular diameter of the Galilean system is about 16–20 arcmin, larger instrument fields (10 arcmin and greater) are advantageous. This will influence the choice of telescope and detector.

B. Galileo Flyby of Asteroid Ida

The key information provided by orbit determination with long arcs (half an asteroid orbit period or more) of ground-based angular star-relative asteroid observations consists of accurate target ephemeris coordinates in three orthogonal directions [5; pp. 31–34]. When coupled with accurate DSN radio tracking of the spacecraft, this enables accurate determination of the spacecraft's time of arrival, which is orthogonal to the spacecraft-target sky-plane, and so is poorly determined by onboard OPNAV imaging. Because of time-of-flight uncertainties, the Galileo Project must schedule a picture mosaic to be sure of capturing a close-up picture of the asteroid. Since most of these pictures will capture only blank sky, there is a definite need to improve the ground-based asteroid observational accuracy so that the near-encounter picture budget can be used for actual observations of the asteroid.

Galileo has already successfully concluded a historic first encounter with asteroid 951 Gaspra on October 29, 1991. Pre-encounter ground-based astrometric observations obtained by the astronomical community and analyzed at JPL³ provided critical Gaspra target-location improvement to enable accurate Galileo spacecraft instrument pointing. Recent star-relative observations of Gaspra from the NOFS 1.55-m sidereal CCD instrument and 0.2-m CCD transit instrument were a major contributor to the success of the encounter navigation, which achieved orbit prediction errors of less than 100 km (similar to the standard errors from the solution covariance matrix). This provided improved instrument pointing accuracy and enabled successful acquisition of several close-up images of Gaspra.

² D. W. Murrow, "Integrated Covariance for the Galilean Satellites from E3," JPL Interoffice Memorandum 314.3-779 (internal document), Jet Propulsion Laboratory, Pasadena, California, January 21, 1988.

³ D. K. Yeomans, P. W. Chodas, M. S. Keesey, W. M. Owen, Jr., and R. N. Wimberley, "Ground-based Ephemeris Development for Asteroid 951 Gaspra," JPL Interoffice Memorandum 314.6-1417 (internal document), Jet Propulsion Laboratory, Pasadena, California, March 24, 1992.

There will also be a close encounter with asteroid 243 Ida on August 28, 1993. An observation program similar to the Gaspra program is being conducted for Ida. All these observations are limited by inaccuracies in the available star catalogs, so that the star-relative observation noise is about 1450 nrad (0.29 arcsec)—a large value when compared with the technology development goal of 25–50 nrad per night. This motivates an effort to obtain a more accurate star catalog, both to support navigation technology development and to improve the Ida target prediction accuracy.

Fortunately, the European Space Agency (ESA) Hipparcos Earth-orbiting observatory has been acquiring star observations for a global star catalog since late November 1989. Recent estimates [6] of expected catalog accuracy, assuming three years of data, predict star positional accuracies better than $10 \text{ nrad} + (10 \text{ nrad/yr}) T$, where T is the time in years past the end of the catalog data span. Catalog density would be about 2.5 stars per square degree, roughly in an even distribution over the celestial sphere. Eventually, of course, a second Hipparcos mission would be required to reduce the effect of the 10 nrad/yr star proper motion errors, but, in any case, the Hipparcos catalog is expected to enable dramatically improved astrometry.

A special Hipparcos input catalog (not based on Hipparcos data) was obtained from ESA, containing approximate coordinates of 50 Hipparcos stars lying near the Ida track on the sky. Further processing⁴ at JPL identified about nine stars within 10 arcmin of the Ida track during the Spring 1992 observing season, including an opportunity to observe Ida simultaneously with two Hipparcos stars for over a week during the post-opposition stationary point.

Observations of Ida relative to several of these Hipparcos stars were obtained by the NOFS 1.55-m telescope (CCD with 11-arcmin field), including observations taken over several days capturing Ida and two Hipparcos stars. The appearance of two stars with Ida during an Ida stationary point was a fortuitous but highly unlikely event, which may allow accurate scale and orientation calibrations for these observations and thus provide accurate angular positions of Ida, possibly with 50 nrad or better accuracy.

One of the actual NOFS CCD observations, taken in early 1992, is shown in Fig. 3. Ida is the faint (magnitude

⁴ W. M. Owen, Jr., "Opportunities for Observing 243 Ida Relative to Hipparcos Stars," JPL Interoffice Memorandum 314.8-819 (internal document), Jet Propulsion Laboratory, Pasadena, California, January 20, 1992.

$m_v = 14.0$) object surrounded by the cursors, and the two Hipparcos stars ($m_v = 8.3$ and 8.8) are the brightest and largest star images in the field. While the large difference in brightness between Ida and the Hipparcos stars may possibly degrade the measurement accuracy, it is important to note that the NOFS 2048×2048 CCD is the first astrometric CCD with wide enough dynamic range to enable these measurements. For example, an 800×800 CCD of the type used by Galileo or the Hubble Space Telescope would have been saturated by more than a factor of 6 under similar circumstances.⁵

Initial ESA plans were for a four-year data-reduction interval between the last Hipparcos observation (which depends on the spacecraft's lifetime) and final catalog release. The Galileo Project has requested preliminary catalog positions for the observed stars (based on actual Hipparcos data), but it is still not known whether it is feasible to obtain these positions prior to the August 1993 Ida encounter.

If early catalog delivery is not possible, then there will be no navigation benefit for the Ida encounter. However, for technology demonstration purposes, catalog delivery after Ida encounter would still provide a useful demonstration, since Ida's heliocentric orbit will be known very accurately after the encounter, so that computed residuals for the Hipparcos-relative observations could be used to verify the data accuracy.

III. Instrument Overview

This section reviews the characteristics of current photographic, Ronchi, and CCD astrometric instruments, reviews their capability to accurately observe the Galilean satellites, and provides the rationale for choosing a CCD instrument for the present Galilean-satellite observing-technology demonstration.

For readers desiring more information about CCD or Ronchi systems, there is an excellent astrometric review by Monet [7]. Readers who wish to accept the outcome of this section (i.e., the choice of the NOFS CCD for a Galilean satellite technology demonstration) can skip to Section IV, Galilean Satellite CCD Observation Techniques, without losing any information required for understanding the rest of this article.

A. Photographic Instruments

Photographic detectors have serious systematic defects (nonlinear response and emulsion shifts), and their quan-

tum efficiency is only a few percent. In practice, the non-linearity leads to magnitude-dependent, position-dependent changes in star image locations; adequate calibration of these effects is very difficult. Nevertheless, photographic techniques enable observation over wide fields; these techniques have been extensively used for Galilean intersatellite observations since the early 1900s, with errors in the best cases as small as ± 250 nrad (0.05 arcsec) [8].

Recent NOFS photographic observations and analyses by Pascu et al. [9] show Galilean intersatellite errors of about 200 nrad for angular separation $S < 100$ arcsec, and of about 550 nrad for a complete set of measurements, including many at larger separations. The latter value is more applicable to the present navigation application, since, as discussed, it is important to observe the entire region out to $S \leq 600$ arcsec.

Pascu et al. attribute this error pattern to scale errors, whose effect is proportional to S . This can be confirmed by acquiring repeated observations of bright star fields, then examining the night-to-night reproducibility of the plate-constant solution residuals. Such results effectively remove all errors that are functions of star position, color, or brightness, and therefore the reproducibility is an underestimate of the actual observational errors.

For example, photographic observations of the star field surrounding 51 Andromedae were acquired and analyzed by Stein and Castelaz⁶ in support of the present JPL technology development effort. This work, performed with the Allegheny Observatory 0.76-m reflector (a different telescope from the Allegheny 0.76-m refractor used for Ronchi-ruling observations), showed night-to-night reproducibility of about 100 nrad. Stein and Castelaz also found that repeated photographic observations of Galilean satellites (all taken within about 1 hour) showed reproducibility of about 150–250 nrad, after being reduced to a common scale value and extrapolated for satellite motion. Observations of asteroid 243 Ida ($m_v = 15$) were unsuccessful because of various difficulties in observing such a faint object with the 0.76-m photographic reflector.

All these results suggest that it might be possible to achieve 100- to 200-nrad photographic single-observation accuracy for the Galilean satellites, if accurate calibrations can be made for instrument scale and systematic

⁵ D. G. Monet, *op. cit.*

⁶ J. W. Stein and M. W. Castelaz, *Acquisition and Data Analysis of Ronchi Ruling and Photographic Data and Crossing Point Measurements of Asteroids*, Final Report on Proposal IC-6-8063, work performed for the Jet Propulsion Laboratory, at the Allegheny Observatory of the University of Pittsburgh, Pittsburgh, Pennsylvania, 1991.

errors. However, navigation system development eventually will require 25-nrad observational accuracy and a capability to observe faint sixteenth-magnitude objects, and this accuracy appears to be well beyond the capabilities of photographic techniques. Therefore, modern photoelectric detectors appear to provide the best opportunity for a technology demonstration that could support long-term development plans and potentially provide better accuracy than with photographic detectors.

B. Ronchi Ruling Instruments

Ronchi ruling devices receive the light from each object (star or solar system body) in a separate photometer. This light is modulated by a moving grating, consisting of precisely ruled alternating opaque and clear strips. The difference in modulation phase between the various observed objects can be analyzed to provide accurate angular differences in the direction of grating motion. Two observers are actively using Ronchi astrometric instruments. The first, A. Buffington [10] (Table Mountain Observatory), has a fixed meridian-transit instrument (0.29-m reflector), which currently can observe only in right ascension and is very sensitive to scattered light from bright sources, such as Jupiter. The second, G. Gatewood [11] (Allegheny Observatory), has a sidereally guided 0.76-m refractor, but this instrument would require major modifications to keep the image of the moving satellites in the photometers during an exposure. Since observations in both right ascension and declination are required, moving bodies must be observed, and major near-term instrument renovations should be avoided, neither of these instruments is suitable for the present target-location applications.

C. CCD Instruments

Finally, CCD instruments potentially can accurately observe the Galilean satellites. These instruments have high quantum efficiency, essentially linear response, and a stable, precise metric. However, as discussed, it is important to have a very large format CCD, so that a good trade-off of field size and pixel size can be achieved. The NOFS 1.55-m instrument, which NOFS recently upgraded with a 2048×2048 chip for their own purposes, is currently unsurpassed in this regard, with an 11-arcmin field and a small, roughly 0.3-arcsec pixel size. Therefore, the technology demonstration for Galileo is being conducted with this instrument.

IV. Galilean Satellite CCD Observation Techniques

This section presents an overview of narrow-field CCD observational techniques, in the context of obtaining accu-

rate ground-based intersatellite observations of the Galilean satellites. Some analysis results will also be presented. As discussed, the 1- σ accuracy goal is about 50 nrad per night.

A. Scale-Determination Techniques

Ground-based instrument scale changes significantly with temperature and other environmental conditions, with nightly changes in fractional scale [8, p. 75] of as much as 10^{-4} . This would induce unacceptable 300-nrad errors for 600-arcsec angular separations.

By scale-change calibration with measured temperatures and coefficients of expansion for the focal plane surface, it may be possible to improve this situation slightly for instruments such as the NOFS 1.55-m telescope, which has flat secondary mirrors. For example, analysis of three seasons of 51-Andromedae star field observations, taken with the Allegheny Observatory 0.76-m refractor and Ronchi instrumentation, showed⁷ that variations in temperature-calibrated fractional scale had a standard error of about 3×10^{-5} . This corresponds to about 90 nrad (0.018 arcsec) over a 600-arcsec field, still not accurate enough for present purposes. Since environmental calibrations do not appear to provide sufficient scale accuracy for the 50-nrad per night accuracy goal, it will be necessary to determine the scale roughly coincident with each astrometric observation.

Traditionally, astronomers have determined the instrument scale value by simultaneously observing two or more stars, whose positions must be available from a star catalog. Then, the scale in arcsec/pixel can be computed, because for each star both the angular position (in arcsec, from the star catalog) and linear position (in pixels, from image centroids) are available. Figure 4 shows a multistar sky-plane observing geometry with three stars. Since this geometry provides a large angular separation in two orthogonal sky-plane directions, the instrument scale is well determined in all directions.

However, narrow-field CCD instruments usually cannot observe any stars without overexposing a bright target such as a fifth-magnitude Galilean satellite (stars with dim asteroids are much easier to observe). Even if the instrument is pointed away from the target so that long exposures bring up stars, then the observed faint stars will usually not have accurate a priori positions.

⁷G. W. Null, "Preliminary Analysis of Allegheny MAP Data (51-Andromedae Star Field)," JPL interoffice memorandum 314.5-1404 (internal document), Jet Propulsion Laboratory, Pasadena, California, February 1, 1990.

NOFS is investigating various combinations of two generic techniques to overcome these narrow-field scale-determination problems.⁸ The first uses catalog densification with another, wider field instrument, and the second observes image motion across the CCD field.

Catalog densification involves observation of numerous faint stars relative to a few stars whose catalog positions are accurately known. This provides a densified catalog with many accurately known stars, which in turn provides three or more well-distributed, accurately known stars in the same CCD field with the target body. Then the scale can be computed. Densification usually requires a wide field to capture sufficient numbers of bright stars.

NOFS is performing catalog densification with the NOFS 0.2-m transit instrument [7, pp. 428 and 432] which clocks out the CCD charge at the sidereal rate. This instrument potentially could achieve 50-nrad accuracy.

Image motion scale-determination techniques take advantage of the fact that, although the angular position of a star or target body may be poorly known, angular motion expressed as the difference of positions at two epochs is usually accurately known. This angular motion can either be from target motion relative to the star background [12] (since target mean motion is usually well known) or apparent motion induced by stopping and starting the telescope drive⁹ (thus making use of the Earth's well-known rotation rate). In either case, the time interval between observations is usually chosen so that the image moves across most of an instrument's field of view. Then the scale is obtained by the ratio of the angular change (in arcsec, from the product of the time interval and angular rate) to the linear change (in pixels, from the brightness centroid at each epoch).

When possible, there should be observations of two or more stars in the field (not necessarily with accurate catalog positions), so that observations taken during the image motion can all be accurately reduced to the same scale. Otherwise, it is necessary to rely on scale stability during this interval.

B. Orientation-Determination Techniques

If two or more catalog stars are in the field of view, then the instrument orientation (i.e., angular orientation about the optical axis) can be determined. Otherwise, the orientation can be obtained by the previously discussed image-motion observations.

⁸ D. G. Monet, *op. cit.*

⁹ *Ibid.*

C. Calibration for Albedo Variations

Analysis [5, p. 21-23] of digitized Voyager Galilean satellite mosaic maps [13] obtained centroid shift versus satellite rotational longitude; the maximum centroid shift caused by albedo variations was about 0.05 satellite radius. The error in extrapolating these results from Voyager's vidicon (strong blue response) to a CCD (strong red response) was found to peak at about 0.01 radius.

These results were obtained using moment centroid algorithms; analysis of simulated ground-based CCD images for the NOFS 1.55-m instrument indicate that when a two-dimensional Gaussian is used for the centroid fit, the maximum effect is reduced by about a factor of three, i.e., to 0.017 radius. This corresponds to about 25 to 45 km, depending on satellite radius, still somewhat larger than the desired 15-km (25-nrad) orbit accuracy. The two-dimensional Gaussian is roughly comparable to centroid fitting functions actually used at NOFS.

Although mosaic maps could be used for albedo-shift calibration, recent analysis¹⁰ has indicated that seams in these maps create unacceptable errors. Work is in progress¹¹ to perform the necessary calibrations with fit-every-pixel techniques applied to the original Voyager satellite images. This effort will be verified by using some Voyager images to predict albedo shifts for other Voyager images. If successful, it should be possible to adequately calibrate the ground-based Galilean satellite images.

V. Summary and Conclusions

This article has described some near-term technology development to support future optical astrometric tracking of laser-emitting spacecraft, target bodies, and stars. Since there currently are no such spacecraft, a good development substitute is to track target bodies relative to each other or to the star background.

A narrow-field CCD observing option has been identified, which can be tested without a major development effort and which potentially can provide significant navigation target-location benefits to the Galileo mission. A cooperative arrangement has been made with USNO Flagstaff Station to obtain and analyze observations with

¹⁰ P. J. Dumont, personal communication, Optical Systems Analysis Group, Jet Propulsion Laboratory, Pasadena, California, April 15, 1992.

¹¹ J. E. Riedel and P. J. Dumont, personal communication, Optical Systems Analysis Group, Jet Propulsion Laboratory, Pasadena, California, April 15, 1992.

a 1.55-m telescope, whose CCD detector provides an 11-arcmin field of view.

The key technology challenge is to demonstrate that it is possible to accurately calibrate the instrument scale, since the narrow field will usually not capture enough catalog stars for traditional scale-determination methods. To test scale-calibration techniques, observations of stars, now in progress, will be analyzed with two alternative scale-determination methods. First, catalog densification with the NOFS 0.2-m CCD transit instrument will provide a local star catalog with enough stars to enable scale determination for the sidereal instrument. Second, image motion across the 11-arcmin CCD field will be used to determine the scale.

Then, after demonstration of adequate scale-determination methods, intersatellite observations will be acquired for Jupiter's Galilean satellites. The 50-nrad per night accuracy goal for these intersatellite measurements appears to be potentially achievable. If acceptable accuracy is achieved, then these observations will be included in the Galilean satellite ephemeris determination.

The ephemeris goal is to reduce the positional standard error (currently 50–100 km) down to about 15 km (about 25 nrad), by combining many observations. This would have an important navigation benefit for the Galileo–Jupiter tour, since an accurate position of the spacecraft relative to the target satellite could be obtained with a much smaller set of onboard-optical pictures than would otherwise be required. This would enable more science pictures to be transmitted to Earth with the currently restricted spacecraft antenna configuration.

Current efforts to obtain NOFS CCD observations of Galileo asteroid target 243 Ida relative to the ESA Hipparcos star catalog were also discussed. A few nights of observations containing 243 Ida and two Hipparcos stars were obtained; these observations potentially could provide very accurate angular positions of this asteroid. If the Hipparcos output catalog is available prior to Galileo's August 1993 encounter with Ida, these data could provide valuable navigation improvements; in any case, it probably will be possible to conduct a technology demonstration to demonstrate the accuracy of Hipparcos-relative observations.

Acknowledgments

The authors thank D. G. Monet of NOFS for useful conversations, for his cooperation as the NOFS technical contact, and for providing the NOFS CCD pictures shown in Figs. 2 and 3. They also gratefully acknowledge useful conversations with J. S. Ulvestad of the Navigation Systems Section.

References

- [1] J. R. Lesh, L. J. Deutsch, and W. J. Weber, "A Plan for the Development and Demonstration of Optical Communications for Deep Space," *TDA Progress Report 42-103*, vol. July–September 1990, Jet Propulsion Laboratory, Pasadena, California, pp. 97–109, November 15, 1990.
- [2] K. Shaik, "A Two-Telescope Receiver Design for Deep Space Optical Communications," *TDA Progress Report 42-101*, vol. January–March 1990, Jet Propulsion Laboratory, Pasadena, California, pp. 114–120, May 15, 1990.
- [3] J. H. Lieske, "Improved Ephemerides of the Galilean Satellites," *Astron. Astrophys.*, vol. 82, pp. 340–348, 1980.
- [4] K. Aksnes, F. Franklin, R. Millis, P. Birch, C. Blanco, S. Catalano, and J. Pironen, "Mutual Phenomena of the Galilean and Saturnian Satellites in 1973 and 1979/1980," *Astron. Journal*, vol. 89, no. 2, pp. 280–288, February 1984.

- [5] G. W. Null, *Deep Space Target Location With Hubble Space Telescope and Hipparcos Data*, JPL Publication 88-4, Jet Propulsion Laboratory, Pasadena, California, February 15, 1988.
- [6] M. A. C. Perryman, "Hipparcos: Revised Mission Overview," *Adv. Space Research*, vol. 11, no. 2, pp. (2)15-(2)23, 1991.
- [7] D. G. Monet, "Recent Advances in Optical Astrometry," *Ann. Rev. Astron. Astrophys.*, G. Burbidge, D. Layzer, and J. G. Phillips, eds., vol. 26, pp. 413-440, 1988.
- [8] D. Pascu, "Astrometric Techniques for the Observation of Planetary Satellites," *Planetary Satellites*, J. A. Burns, ed., Tucson: University of Arizona Press, 1977.
- [9] D. Pascu, C. A. Adler, and J. F. Bloomfield, "An Analysis of Photographic Astrometric Observations of the Galilean Moons: USNO Refractor, 1986-1990," *Bull. Amer. Astron. Soc.*, vol. 23, no. 3, p. 1255, 1991.
- [10] A. Buffington and M. R. Geller, "A Photoelectric Astrometric Telescope Using a Ronchi Ruling," *Publ. Astron. Soc. Pacific*, vol. 102, pp. 200-211, February 1990.
- [11] G. D. Gatewood, "The Multichannel Astrometric Photometer and Atmospheric Limitations in the Measurement of Relative Positions," *Astron. Journal*, vol. 94, no. 1, pp. 213-224, 1987.
- [12] R. L. Duncombe, W. H. Jefferys, P. J. Shelus, P. D. Hemenway, and G. F. Benedict, "Astrometry Using the Hubble Space Telescope Fine Guidance Sensors," *Adv. Space Research*, vol. 11, no. 2, pp. (2)87-2(96), 1991.
- [13] T. V. Johnson, L. A. Soderblom, J. A. Mosher, G. E. Danielson, A. F. Cook, and P. Kupperman, "Global Multispectral Mosaics of the Icy Galilean Satellites," *Journal of Geophysical Research*, vol. 88, no. B7, pp. 5789-5805, July 10, 1983.

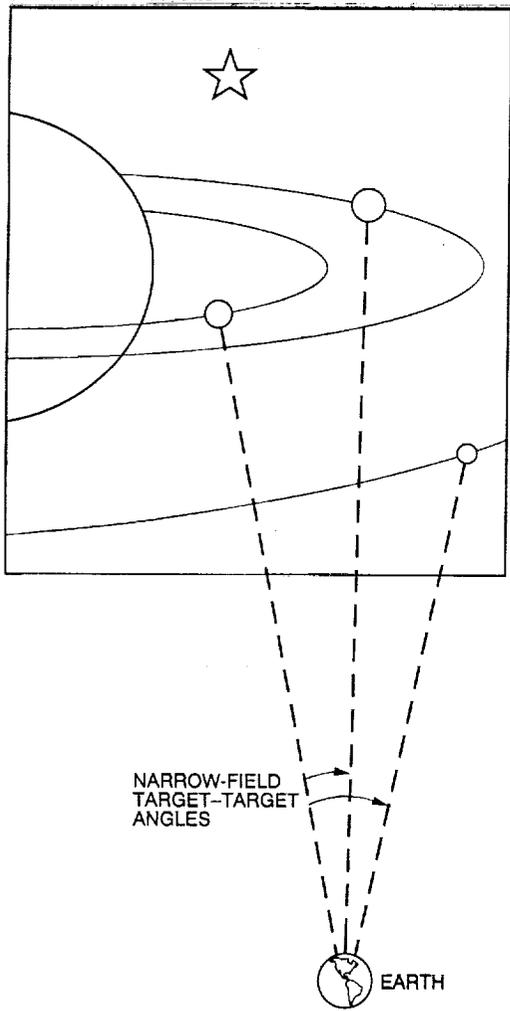


Fig. 1. Galilean satellite observing geometry.

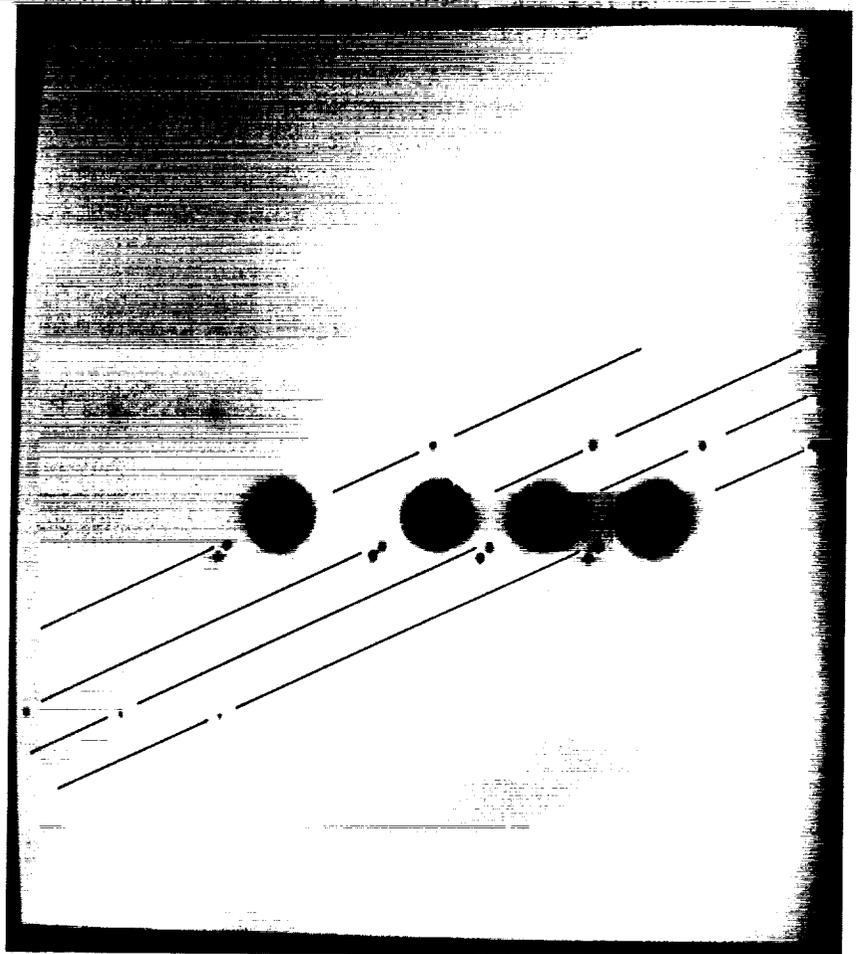


Fig. 2. NOFS multiframe CCD picture of Jupiter's Galilean satellites (courtesy of D. G. Monet, NOFS).

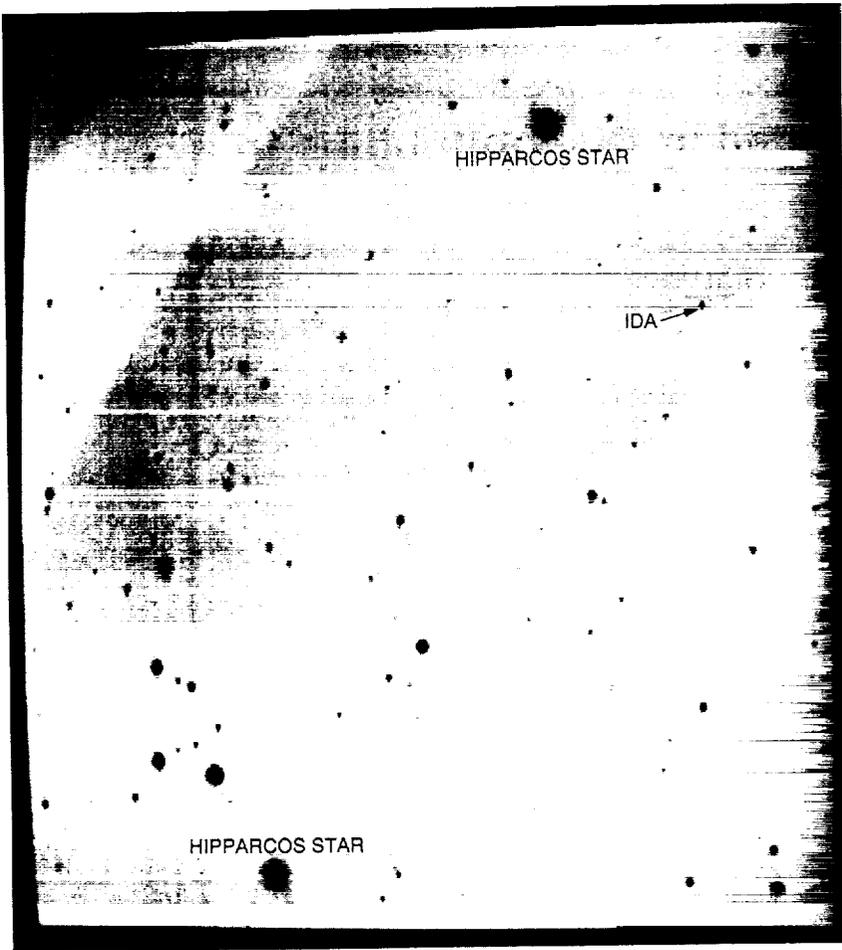


Fig. 3. NOFS CCD picture of asteroid 243 Ida with two Hipparcos stars (courtesy of D. G. Monet, NOFS).

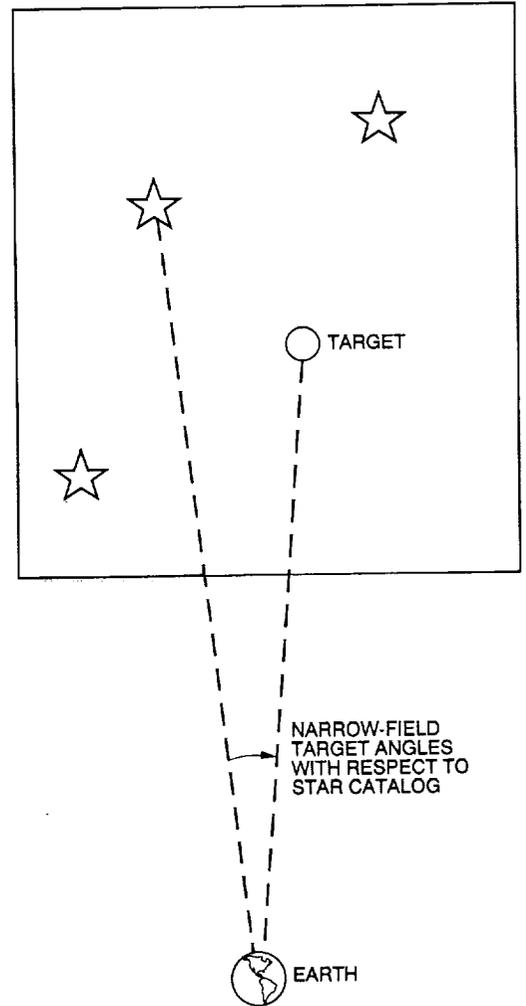


Fig. 4. Target-body observation relative to star background.

510-327-19423
 128443
 P. 23

A Method for Modeling Discontinuities in a Microwave Coaxial Transmission Line

T. Y. Otoshi

Ground Antennas and Facilities Engineering Section

This article presents a method for modeling discontinuities in a coaxial transmission line. The methodology involves the use of a nonlinear least-squares fit program to optimize the fit between theoretical data (from the model) and experimental data. When this method was applied to modeling discontinuities in a slightly damaged Galileo spacecraft S-band (2.295-GHz) antenna cable, excellent agreement between theory and experiment was obtained over a frequency range of 1.70–2.85 GHz. The same technique can be applied for diagnostics and locating unknown discontinuities in other types of microwave transmission lines, such as rectangular, circular, and beam waveguides.

I. Introduction

The Galileo spacecraft, launched on October 18, 1989, is currently on its interplanetary journey to encounter Jupiter in 1995. One of the important experiments to be performed with this spacecraft in 1993 is the gravitational wave experiment to support Einstein's General Theory of Relativity.

Prior to launch, Galileo underwent environmental testing in the Space Simulator at JPL. After these tests, it was discovered that the S-band output power had dropped about 0.2 dB at the transmit frequency of 2.295 GHz. More alarming were the radical changes observed in the subsequent measured insertion loss versus frequency characteristics of the S-band antenna cable. Instead of the small peak-to-peak sinusoidal variations seen previously on the pre-environmental test frequency response curves, numerous nonperiodic humps and valleys (of unusual amplitudes) were seen on the post-environmental test curves.

This author (consulting on an emergency basis) diagnosed the problem and provided a satisfactory mathematical model of the slightly damaged cable. This model depicted the Galileo cable as having developed crimps at two different cable clamp locations, and reasonable agreement between theory and experiment was obtained. Further analysis indicated that if the magnitudes of discontinuities doubled during Galileo's interplanetary journey to Jupiter, an additional 0.2-dB loss increase would occur at the 2.295-GHz transmit frequency, but this worst-case expected loss increase (0.4-dB total degradation) would still be acceptable.

Based upon the author's satisfactory explanation of the phenomenon, the Galileo Project decided that the S-band cable did not have to be replaced, and the spacecraft was shipped to Cape Canaveral on schedule. The decision not to replace the cable was partially based upon the fact that X-band (8.415 GHz) was the prime data channel for ra-

dio science experiments and S-band was less important. Now that S-band has become the prime data channel due to problems with the X-band high-gain antenna, the expected S-band antenna performance characteristics should be reanalyzed.

The purpose of this article is to transfer the technology gained from the successful modeling of the discontinuities in the Galileo S-band antenna cable. This technique can be applied to obtaining models of discontinuities in other types of transmission lines, including rectangular, circular, and beam waveguides. Once a good analytical model is developed, it can then be used for diagnostic purposes or for analytical studies of worst-case situations. For example, models of transmission lines with discontinuities can be used to determine the effects of multiple reflections that are known to degrade the frequency stability and noise temperature performance of a receiving system.

The writing of this article was also motivated by the desire to provide radio scientists and gravitational wave experimenters with an accurate model of the damaged S-band antenna cable on Galileo. Such a model might assist in the error analyses of radio science data, should the need arise in the near future. Discontinuities in a cable give rise to multiple reflections that affect Doppler phase, frequency, and group delay stability. Cable movement and cable temperature changes can cause the phases of individual reflected waves to change, thereby affecting the overall phase and amplitude of the output signal. Continued flexing of the cable during Galileo's interplanetary journey to Jupiter can also make the magnitudes of the individual reflections worse. An accurate model can help to establish error bounds (on radio science data) associated with worst-case signal level and frequency stability degradation situations.

In the following sections, the methodology is described and then demonstrated by applying it to the Galileo cable problem for which experimental data were already available. Comparisons are given showing results obtained with a trial-and-error method versus the proposed new method.

II. Previous Models of the Galileo Cable

As was described in a 1989 report,¹ Model 1 was the first trial model that was derived. It consisted of shunt susceptances (representing discontinuities) separated by different lengths of coaxial line. The first necessary step

¹ T. Y. Otoshi, "Galileo Cable Study Report," JPL Interoffice Memorandum 3328-89-0108, (internal document), Jet Propulsion Laboratory, Pasadena, California, May 17, 1989.

of the modeling procedure was to obtain time domain plots from measured *S*-parameters over a frequency range centered at the particular frequency of interest. The *S*-parameters for the Galileo cable were obtained with the HP 8510B automatic network analyzer. Then the shunt susceptance magnitudes² and approximate locations of discontinuities in the cable were estimated from the return loss-time domain plots.

When the unaltered values (extracted directly from return loss-time domain plots) were used in Model 1, the agreement between theoretical and experimental data was poor. The parameters for the subsequently developed Model 2 were almost the same as those of Model 1, except that the line lengths between discontinuities were adjusted slightly to cause all the individual reflection coefficients to add up in phase at the input port at 2310 MHz. As can be seen in Fig. 1, Model 2 provided satisfactory agreement between theory and experimental data. Figure 2 shows the equivalent circuit for Model 2. From this equivalent circuit, it was determined that a total of four discontinuities existed near the connector regions and cable clamp locations (see Fig. 3).

The agreement between theory and experiment for Model 2 was considered to be satisfactory (even very good) in 1989. However, it was known then that the model was incomplete because at least one more discontinuity had to be added near one of the connectors. This was known because when all discontinuities were removed, except those representing the connectors, the periodicity of the insertion loss versus frequency curve was wrong, and therefore the agreement between the calculated and the pre-environmental cable test data was not good. It appeared that adding one more discontinuity in the connector region closer to the end of the cable would have resulted in a better model. However, further attempts to develop improved models after adding an additional discontinuity and then readjusting the line lengths (by trial and error) for a good fit proved to be excessively time-consuming and fruitless. It was clear that some type of least-squares fitting program was required to obtain a better model. The effort to obtain a better model was temporarily abandoned.

III. Computer Program

Although not tasked to find an improved model for the Galileo cable, the author continued working on the problem on a low-priority basis because of his desire to learn how to perform curve fits between mathematical models

² See Eqs. (A-3) and (A-4) in Appendix A.

and experimental data. It is common practice just to fit experimental data with polynomial curves because a polynomial curve-fitting process is easy to perform and useful for displaying trends and calculating intermediate values. What is not generally known, as has been pointed out by C. Lawson of JPL, is that it is almost as easy to perform curve fits between experimental data and the theoretical results from any mathematical model. If the physical phenomenon associated with the experimental data can be modeled mathematically (no matter how involved and complex), then curve fitting can be done between theoretical and experimental data through the use of a nonlinear least-squares fit (NLSF) program. Variance, correlation coefficients, and standard deviations can also be easily computed from the residuals of the nonlinear curve-fitting process.

A linear model is defined here as one whose coefficients a_i can be expressed explicitly in the form

$$y(x) = a_0 + a_1 f_1(x) + a_2 f_2(x) + \dots + a_i f_i(x) + \dots + a_n f_n(x) \quad (1)$$

Polynomials are a subset of the general linear form given by the above Eq. (1). A nonlinear model differs in that the coefficients to be solved for (best-fitted) can be expressed within any of the expressions for $f_i(x)$ or in almost any mathematical form. As long as the program steps can be written to calculate values of $y(x)$ for input values of x , then a nonlinear least-squares fit can be performed to solve for the unknown coefficients.

With Lawson's assistance, a subroutine computer program was written to find a best-fit model of the Galileo cable by using the existing International Mathematical Scientific Library (IMSL) NLSF computer program already available on the UNIVAC at JPL in 1989. The problem with that program was that it did not allow the user to specify bounds on the parameters to be adjusted and best-fitted. Occasionally the IMSL program gave best-fit parameter results that were not physically realizable.

After further consultation with Lawson, it became clear that JPL needed an improved NLSF program that could be used for various types of ongoing modeling problems at JPL. To fulfill this immediate need, Lawson performed an extensive search of available NLSF programs in industry and academic institutions. As a result, he recommended the use of a particular NLSF program³ that was later de-

³ This program was a later version of the original release of NL2SOL by David Gay and Linda Kaufman, now at AT&T Bell Laboratories.

scribed in another JPL publication [1]. This program was unlike most available least-squares fit programs in that it allowed the user to specify bounds on the parameters to be best-fitted. This program has the additional advantages of being public domain software and can be run on a personal computer. The input data required for this program are (1) the measured data; (2) the theoretical values for the mathematical model calculated from a subroutine; and (3) estimates of the nominal, upper, and lower bound values of the parameters to be best-fitted.

The Fortran subroutine ultimately developed for the Galileo cable-modeling study calculated the overall S -parameters for a cable that had seven or more discontinuities separated from each other by lengths of coaxial line whose attenuation constant and relative dielectric constant could be specified. Program steps were written to calculate the S -parameters of a basic network consisting of a shunt discontinuity and a length of lossy line. The equations for two types of shunt discontinuities used in this study are given in Appendix A. Another subroutine was written to compute the overall S -parameters of two 2-port networks that were cascaded. When the overall S -parameters of two basic networks are cascaded, the cascaded network becomes the equivalent 2-port that is cascaded with the next basic network. The overall S -parameters are again calculated and stored. This procedure was repeated until all seven basic networks were cascaded.

The insertion loss in decibels of the final overall network was then calculated from

$$IL = 20 \log_{10} |S_{21}| \quad (2)$$

where S_{21} is the S -parameter for the transmission coefficient of the overall cascaded network when terminated in a nonreflecting load [2]. The insertion loss is computed at each frequency to form the theoretical data set for the model. Then comparisons are made with the experimental insertion loss data file which was read into the program. The parameters to be adjusted were (1) the discontinuity magnitudes (in terms of shunt susceptance or capacitance values) and (2) the line lengths separating the discontinuities. The program finds the parameter values (within the specified bounds) that give the best fit (using a least-squares convergence criterion) between theoretical and experimental values. Even though the distances between the discontinuities were allowed to be adjusted within specified bounds, the program was written so that the resulting overall length of the cable for the model had to be equal to the actual physical length of the cable.

IV. Improved Models

After the development of Model 2, it was known that additional discontinuities in the cables existed near the connector interfaces. Examinations of detailed drawings of the connector regions revealed locations of potential discontinuities that were not taken into account in Model 2. Figure 4 shows potential discontinuities as being the back ends (or sharp edges) of the rigid Kynar sleeves located about 2.9 inches from the faces of the cable connectors. Bending of the cable at these points can cause crimps or make permanent deep creases on the outer diameter. Another type of potential discontinuity occurs at the connecting regions (Fig. 5) where there are changes in diameter dimensions and dielectric materials within the coaxial transmission line.

When the NLSF program (described in Section III) became available, new models were developed that included the above-described discontinuities. Two likely equivalent circuits for the cable discontinuities are (1) shunt capacitances whose susceptance values are functions of frequencies and (2) capacitive shunt susceptances whose values do not change over the frequency of interest. The models corresponding to these two types of discontinuities are denoted as Models 3 and 4 and are described below.

A. Model 3

Model 3 represents the cable discontinuities as constant shunt capacitances. The constant shunt capacitance might occur in practice when the reduced outer diameter of the cable is squashed over some physical length (e.g., the cable clamp width). For this type of discontinuity, the magnitude of S_{11} of the individual discontinuity changes with frequency. For this model with nominal values, the final best-fit parameters determined by the NLSF program are shown in Fig. 6. Note that the locations of discontinuities of the model occur very close to actual locations of the cable clamps, the edges of the Kynar sleeve, and the connector discontinuities.

Figure 7 shows the comparison between theoretical and measured insertion losses. It can be seen that the agreement is significantly better than that of the Model 2 fitted curve shown in Fig. 1. To obtain such a good fit, it was necessary to use accurate values of cable attenuation due to line losses between the shunt discontinuities. Appendix B describes the procedure that was used to accurately determine cable attenuation (minus connectors). Although Model 3 appears to be an excellent model, it was found to be incorrect. As may be seen in Fig. 8, when all the discontinuities are removed except the two outer connector discontinuities, the agreement between theory and the

pre-environmental test data becomes progressively worse at the higher frequencies.

B. Model 4

Model 4 is based upon representing the discontinuities as constant susceptance values over the frequency range of interest. This type of discontinuity could be a deep crease, or crimp, on the outer diameter of the cable. Such a discontinuity in practice can be created by bending the cable against the edge of a cable clamp or the edge of a Kynar sleeve. This type of discontinuity can be represented as two shunt capacitances separated by a short distance (less than 0.005 wavelength, or a single shunt susceptance, as discussed in Appendix C). The overall equivalent circuit and locations of the discontinuities along the cable are shown in Fig. 9. Model 4 is the result of best-fitting 15 parameters (seven discontinuities and eight line lengths) by using 101 frequency points for a frequency range of 1.7 to 2.85 GHz. The Model 4 discontinuity locations, shown in Fig. 9, correspond very closely with the actual locations of cable clamps, the edges of the Kynar sleeve, and internal discontinuities of the connectors. The locations of the modeled discontinuities are estimated to be within ± 0.5 in. of the actual cable discontinuities. As with Model 3, the line losses between discontinuities for Model 4 were properly accounted for.

It can be seen in Fig. 10 that the theoretical values for this model agree well (within ± 0.02 dB) with experimental data. Good agreement between theoretical and experimental return losses was also obtained, as can be seen in Figs. 11–12. Figure 13 shows that when all internal discontinuities except the two outer discontinuities are removed, the agreement between the model and pre-environmental test data is still excellent.

To ensure that Model 4 is, in all respects, the correct model of the cable, the phases of S_{21} were considered, as well as the magnitudes. However, attempts to fit both magnitudes and phases to experimental data were unsuccessful. Also, the attempt to do similar least-squares fitting to the S_{11} and S_{22} experimental data proved unsuccessful. It was later revealed that special adapters had been used to measure the reflection coefficients S_{11} and S_{22} . When the effects of the adapters were gated out,⁴ the measurement planes might not have been properly referred back to the connector interface planes of the cable.⁵

⁴ See Hewlett Packard 8510B Operating Manual for a discussion of gating techniques for the purposes of removing external discontinuity effects from the measured S -parameter data.

⁵ W. Folwell, personal communication, Spacecraft Telecommunications Equipment Section, Jet Propulsion Laboratory, Pasadena, California, 1990.

It is difficult to best-fit phase data if the measurement and model reference planes do not coincide.

Due to the difficulty of obtaining a good fit between the theory and the model, based on both S -parameter amplitude and phase data, another method was used to verify the model. This method compared experimental and theoretical time domain plots. Time domain plots require both magnitude and phase information over a wide frequency range. Figures 14-16 show time domain plots for S_{21} , S_{11} , and S_{22} data, respectively. Good agreement was obtained for S_{21} , but for the S_{11} and S_{22} time domain plots, the reference planes had to be shifted by an amount approximately equal to the lengths of the special adapters used in the measurement. The necessity to shift reference planes for S_{11} and S_{22} is consistent with the possibility that the measurement plane might not have been correctly referenced back to the cable connector interface planes.

Despite the described difficulties with the reference plane problem associated with the S_{11} and S_{22} time domain plots, good agreement was obtained for the S_{21} time domain plot. One can conclude from the data presented that Model 4 is an excellent representation of the Galileo S-band cable after the environmental tests.

The sharp edges of the cable clamps and the Kynar sleeves are the probable causes of the discontinuities on the outer diameter of the cable. As was stated in the previously cited report,⁶ the cable clamps should be redesigned

⁶ T. Y. Otoshi, op. cit.

so that crimping or creasing will not occur when the cable is bent against the clamps. It was also stated that only a 0.02-in. reduction in outer diameter of the Galileo cable could produce discontinuities of the magnitudes presented in the models studied. For future spacecraft cables, the edges of the currently rigid Kynar sleeves should also be redesigned and be made flexible.

V. Conclusions

A model has been found that gives excellent agreement between theory and experiment for the Galileo spacecraft S-band antenna cable. The current model now has seven discontinuities (including connectors) instead of four obtained for the previous model.

The excellent results described would not have been possible without the use of the NLSF program. This program was easy to use, and formal documentation⁷ is now available. Engineers are encouraged to use this program for their modeling work.

The method presented here was demonstrated for a coaxial cable, but the technique can be extended to the modeling of discontinuities in other types of transmission lines, such as rectangular, circular, and beam waveguides (with a shroud).

⁷ Math77, Release 4.0, *A Library of Mathematical Subprograms for Fortran 77*, JPL D-1341, Rev. C (internal document), Jet Propulsion Laboratory, Pasadena, California. May 1992.

Acknowledgments

Dr. C. Lawson, at that time supervisor of the Applied Mathematics Group at JPL, graciously supplied the public domain NLSF software and provided assistance with the coding procedures. Scott Stewart, a contractor from PRC in Pasadena, California, programmed the HP 8510B Network Analyzer computer to read theoretical S -parameter data into data files for the purpose of generating theoretical time domain plots. It should be pointed out that this capability does not exist using software available from Hewlett Packard. Technical discussions with Kent Kellogg and Phil Stanton of JPL were helpful.

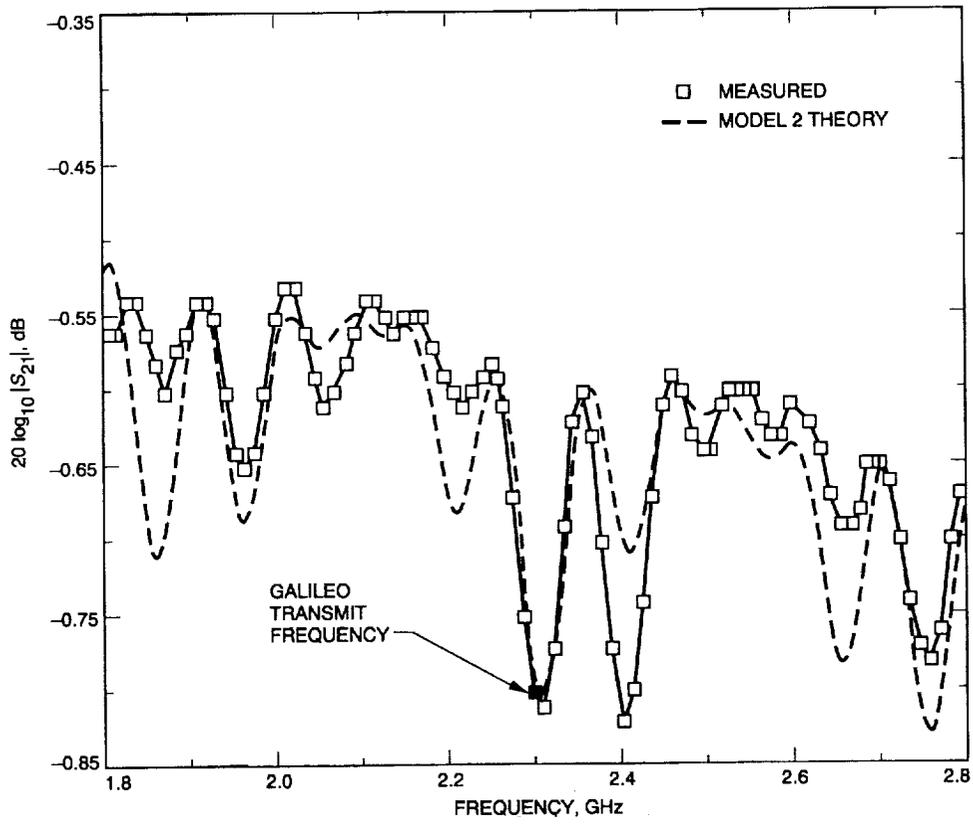
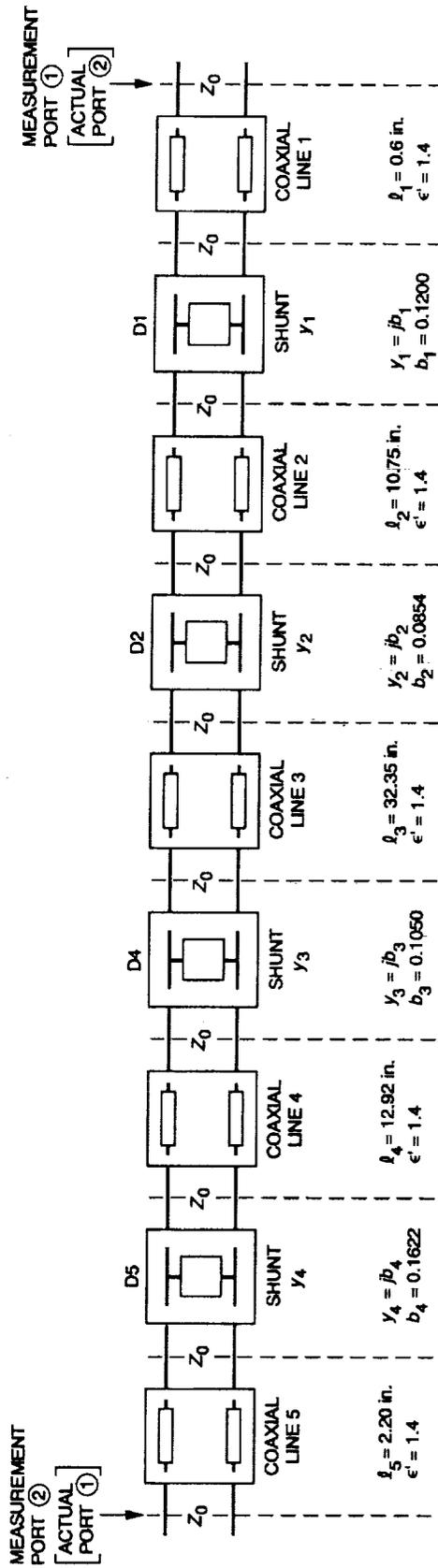
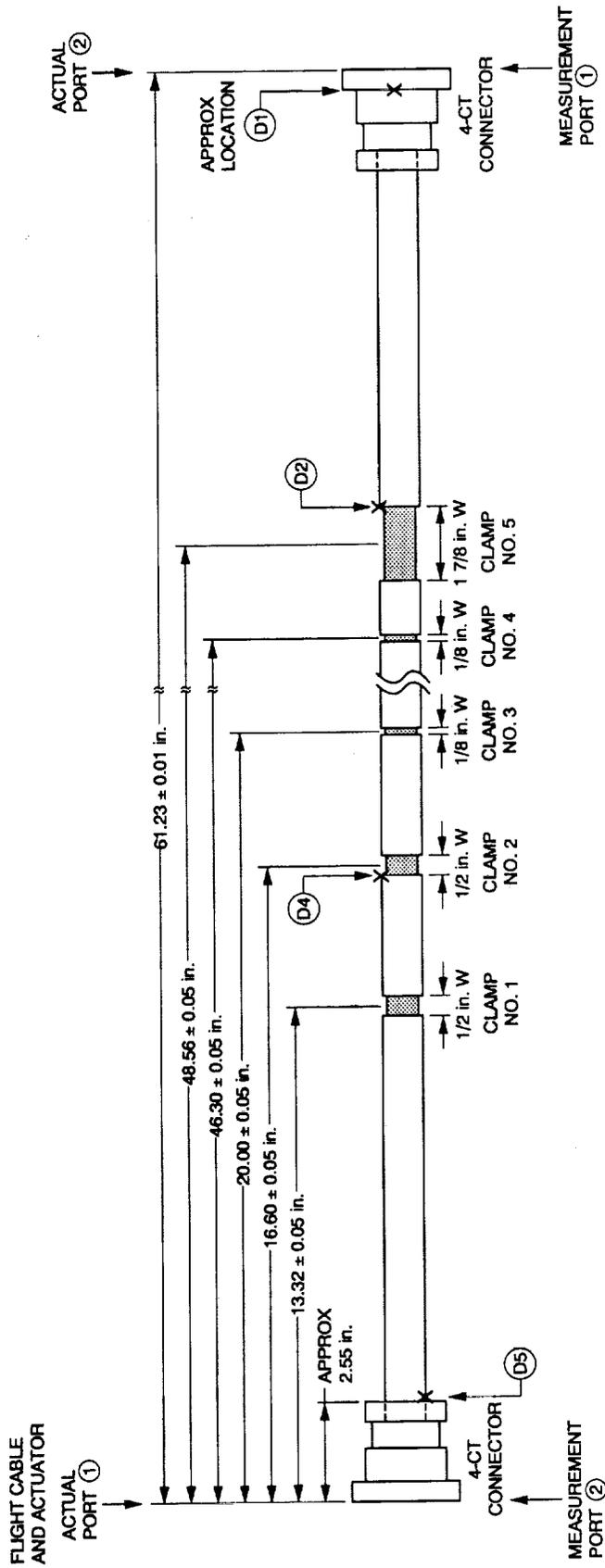


Fig. 1. Comparison of experimental and Model 2 theoretical data.



NOTES: INNER DIAMETER OF OUTER CONDUCTOR IS 0.241 in.
OUTER DIAMETER OF INNER CONDUCTOR IS 0.090 in.
LINE LOSS IS FREQUENCY DEPENDENT AND ACCOUNTED FOR AT EACH FREQUENCY

Fig. 2. Model 2 equivalent circuit of the Galleo S-band cable.



NOTE: D1, D2, D4, AND D5 CORRESPOND TO PEAK RESPONSES ON THE TIME DOMAIN PLOTS - THE LOCATIONS ARE APPROXIMATE BUT ARE ACCURATE TO ± 0.5 in. FOR D2 AND D4

Fig. 3. Model 2 locations of discontinuities relative to locations of cable connectors and clamps.

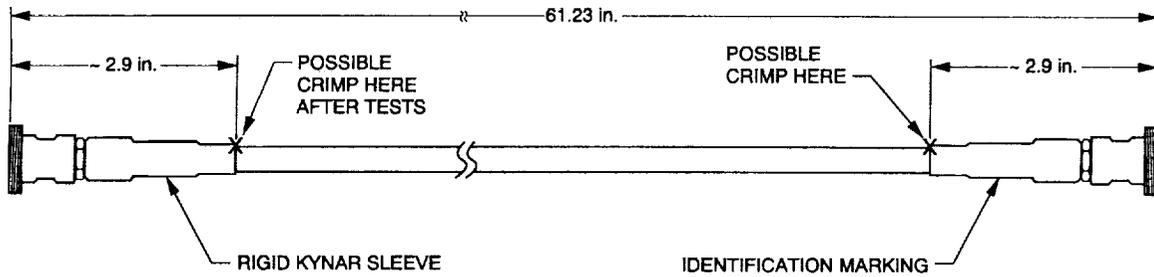


Fig. 4. Galileo S-band cable outer dimension detail.

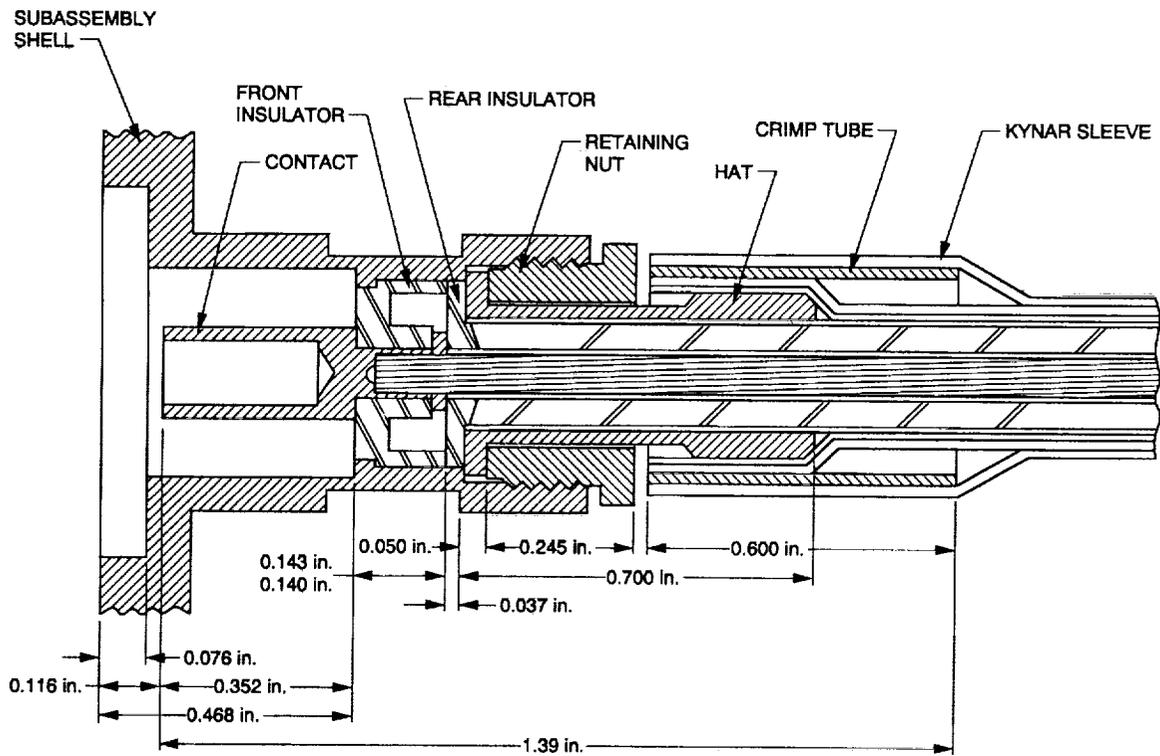


Fig. 5. Connector details.

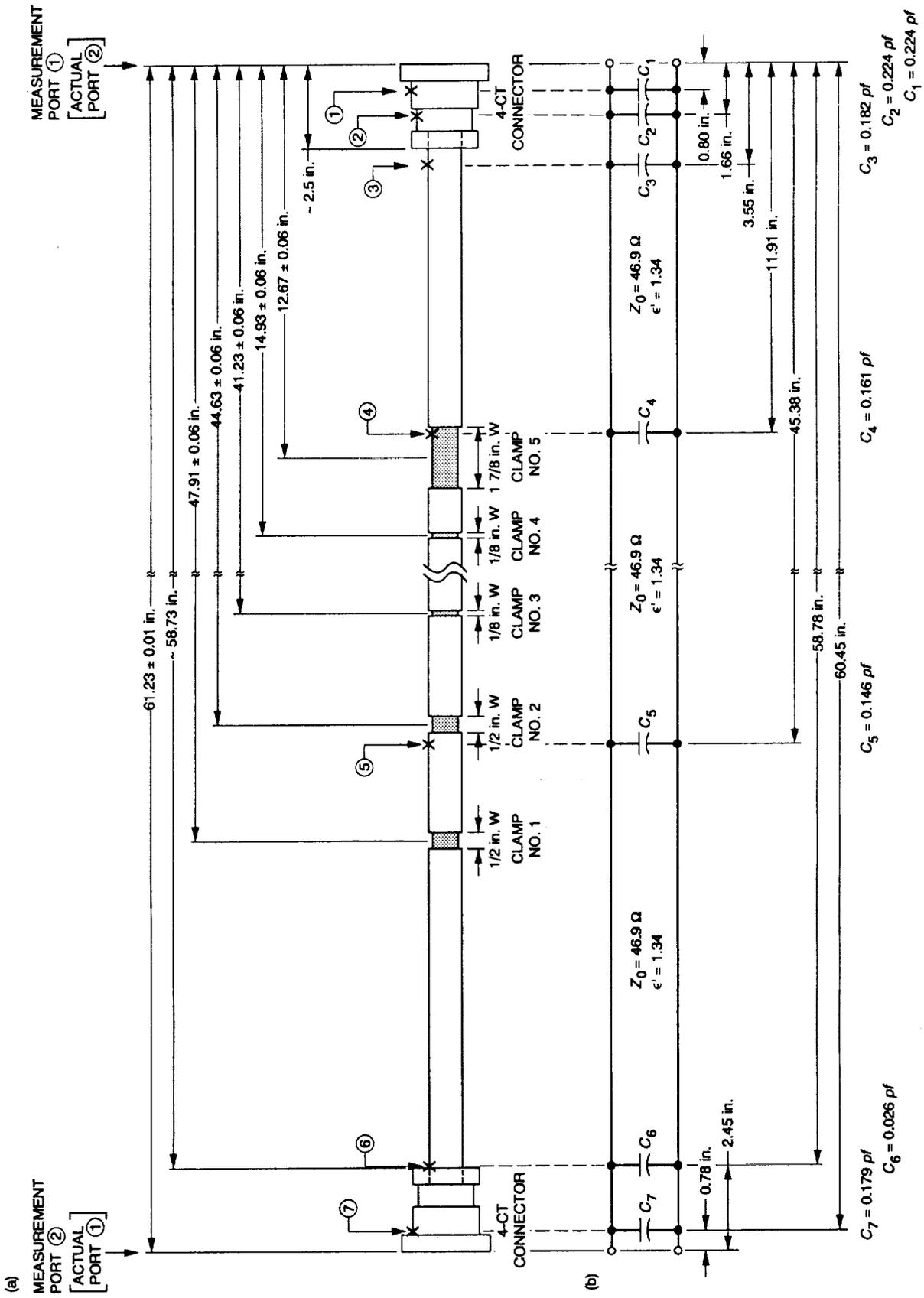


Fig. 6. Model 3: (a) location of discontinuities relative to connectors and cable clamps and (b) equivalent circuit.

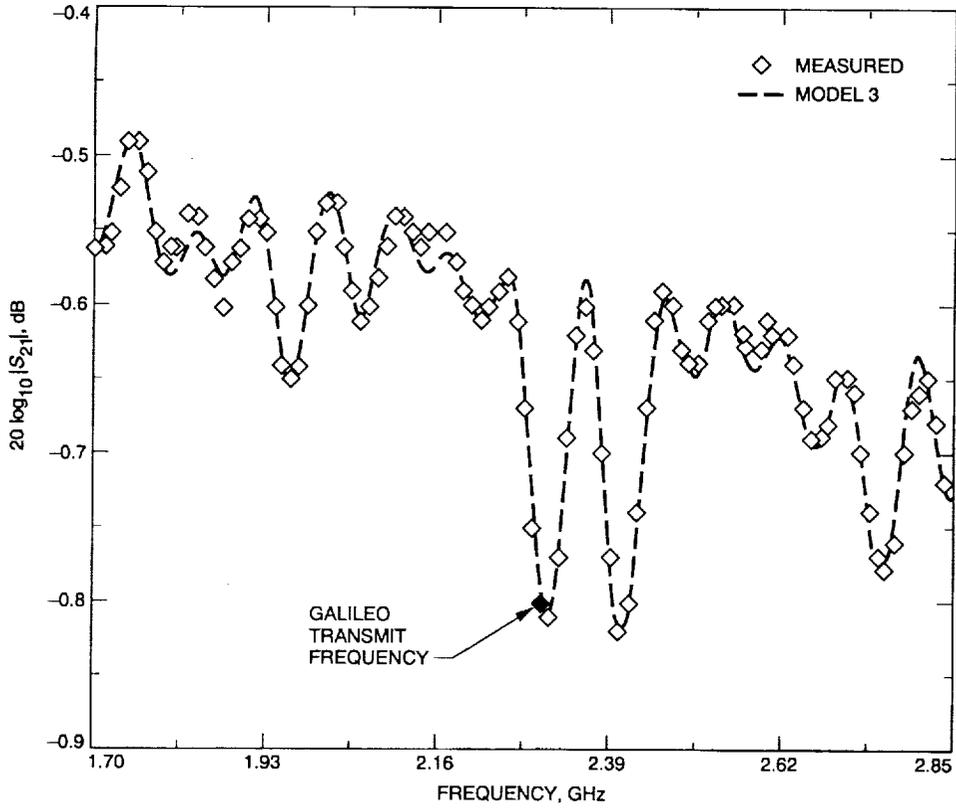


Fig. 7. Comparison of theoretical and measured insertion losses for Model 3, assuming shunt capacitance values that are constant over the frequency range of interest.

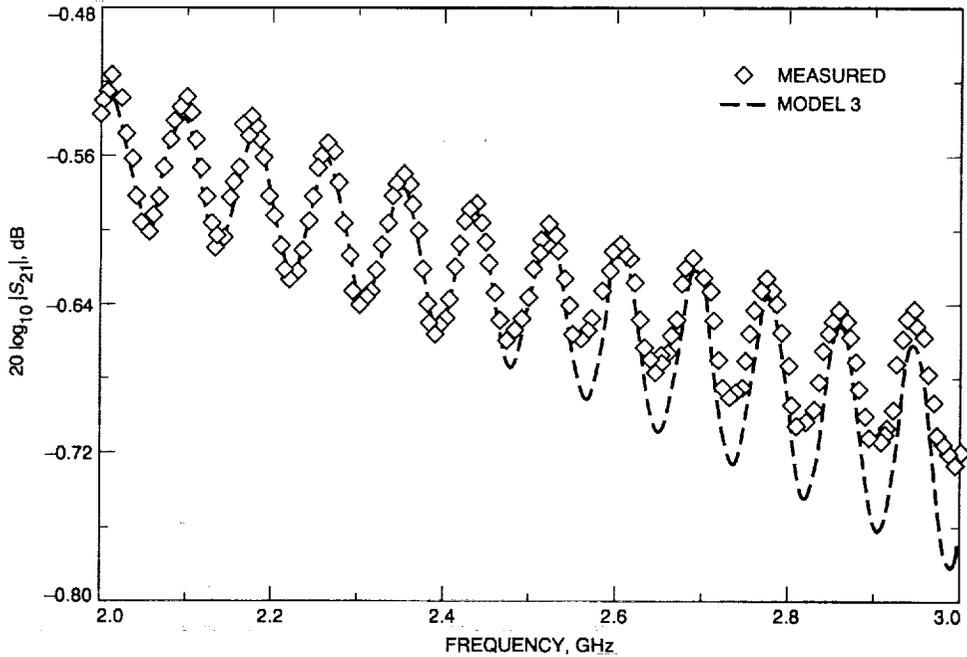


Fig. 8. Comparison of theoretical and measured insertion losses for Model 3 for the pre-environmental cable condition.

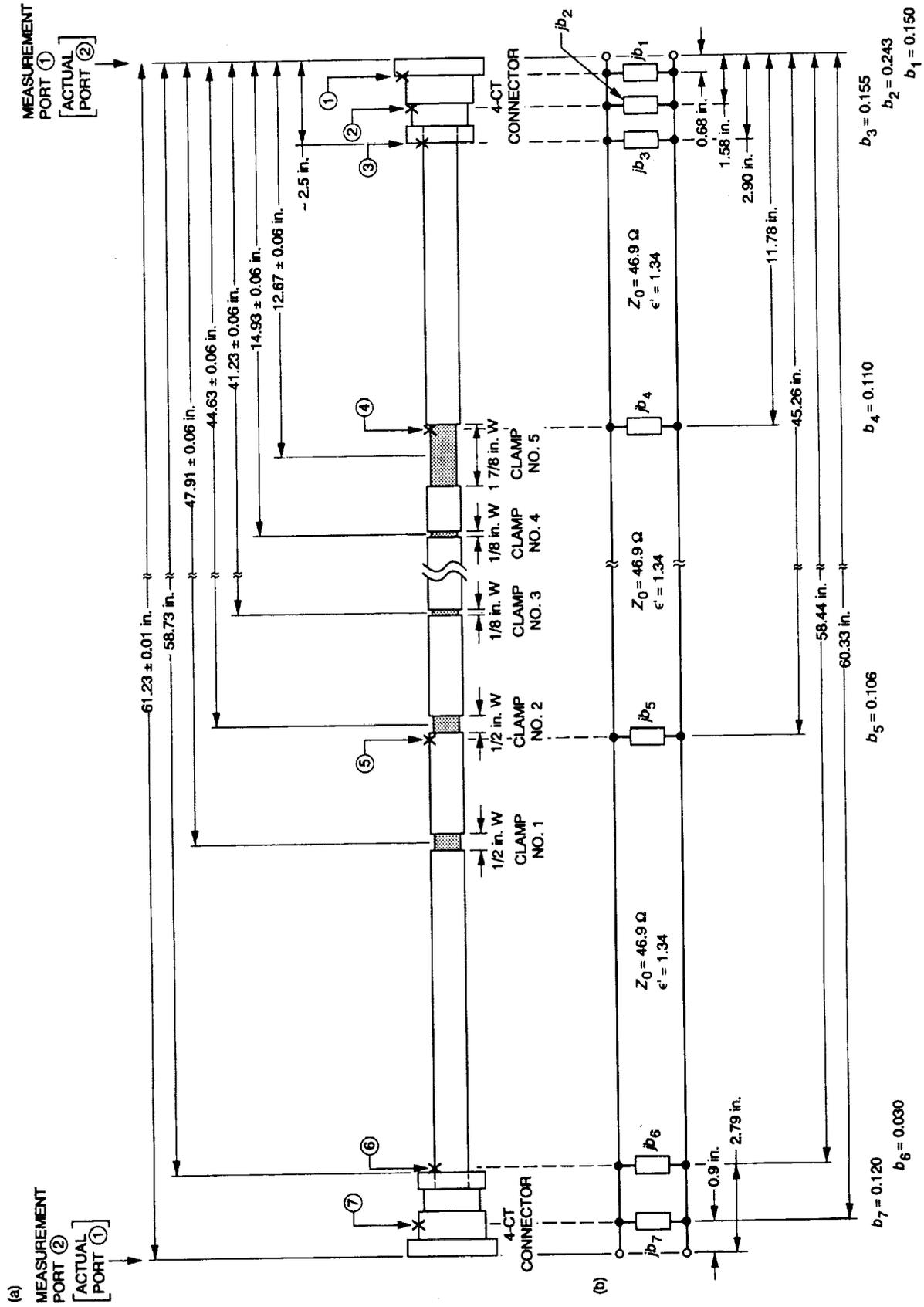


Fig. 9. Model 4: (a) location of discontinuities relative to connectors and cable clamps and (b) equivalent circuit.

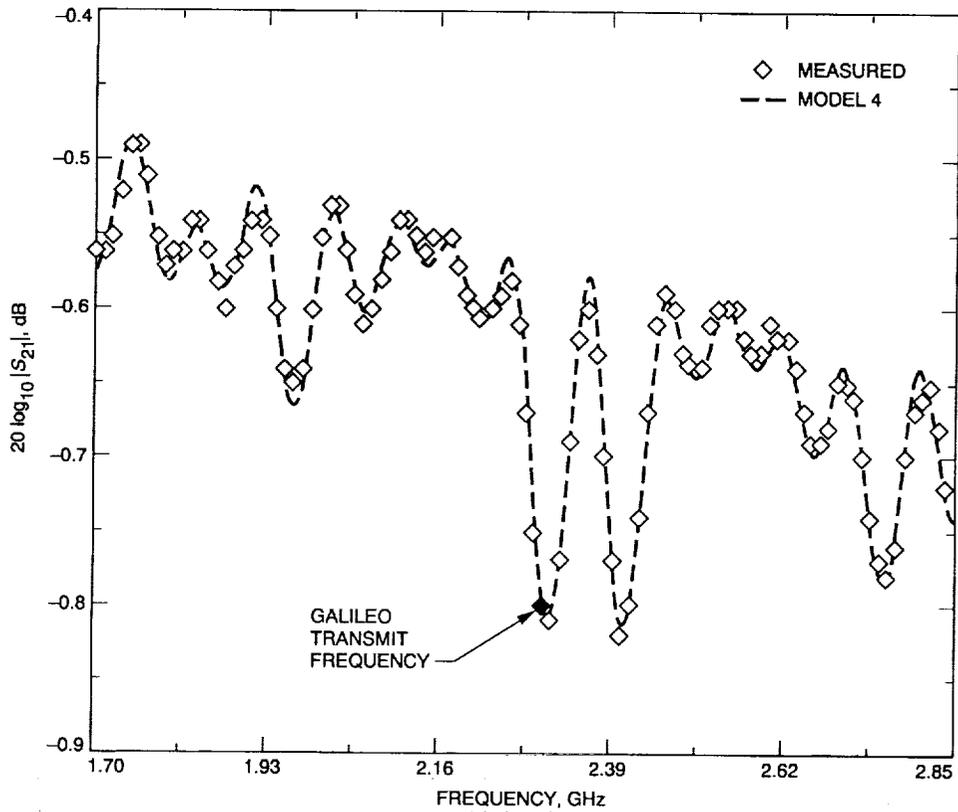


Fig. 10. Comparison of theoretical and measured insertion losses for Model 4, assuming capacitive shunt susceptance values that are constant over the frequency range of interest.

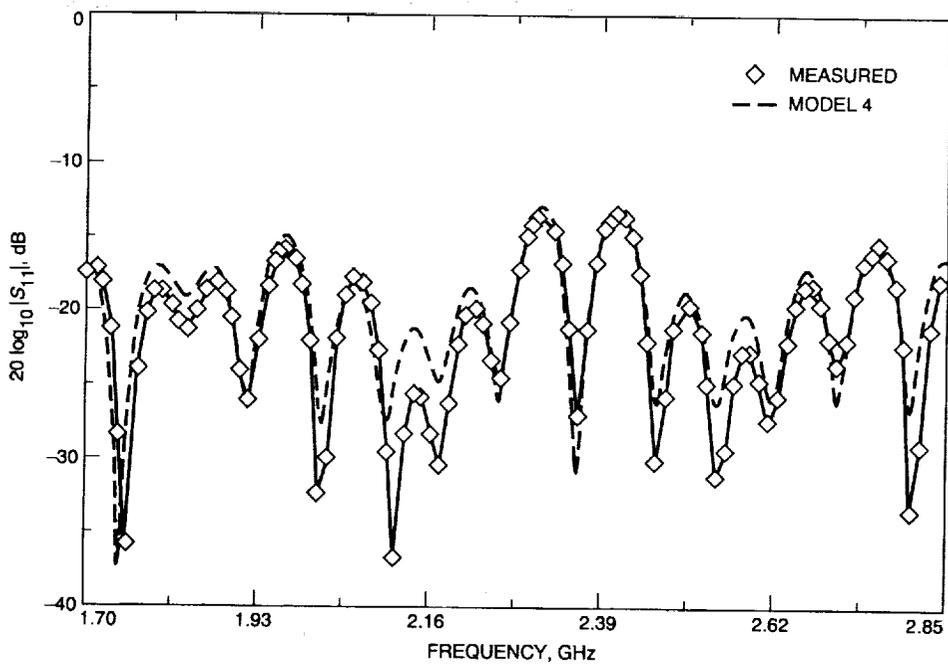


Fig. 11. Comparison of theoretical and measured return losses for Model 4, as seen looking into port 1.

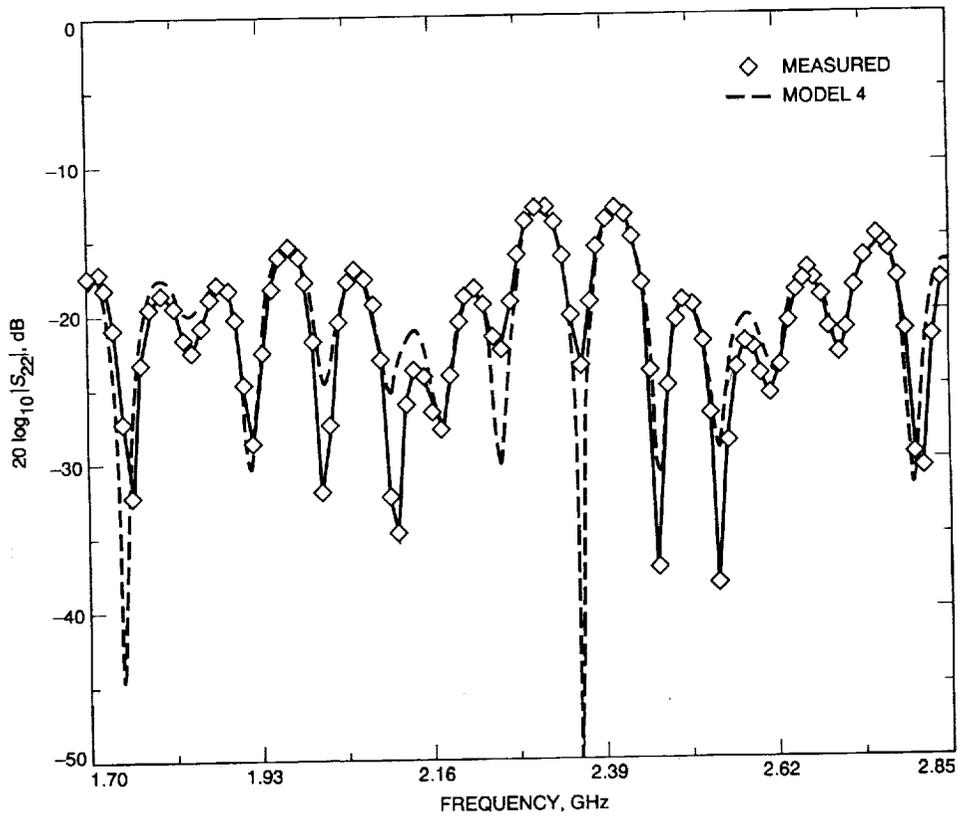


Fig. 12. Comparison of theoretical and measured return losses for Model 4, as seen looking into port 2.

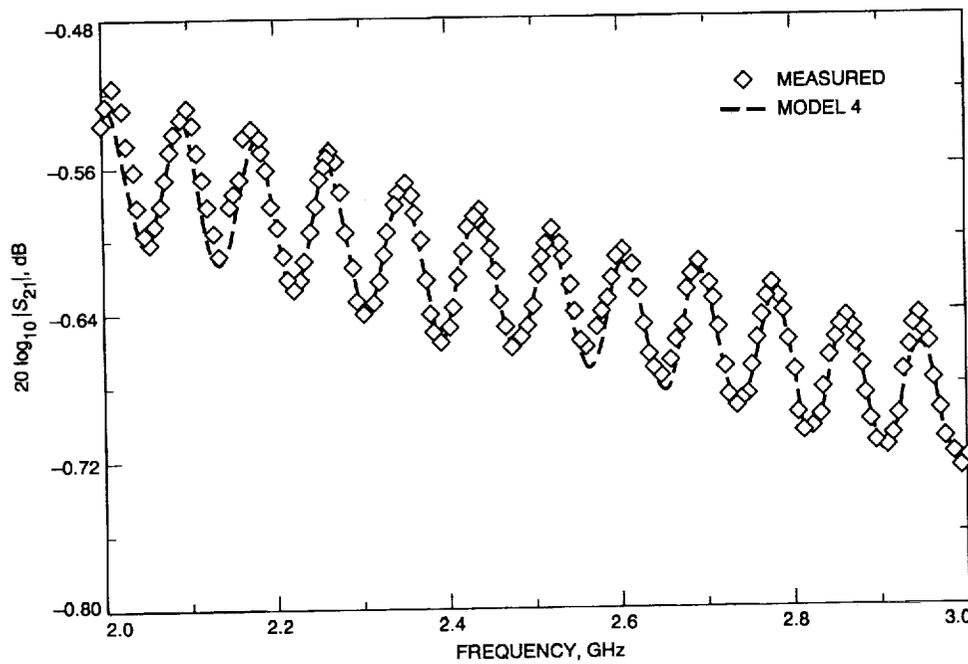


Fig. 13. Comparison of theoretical and measured Insertion losses for Model 4 for the pre-environmental cable condition.

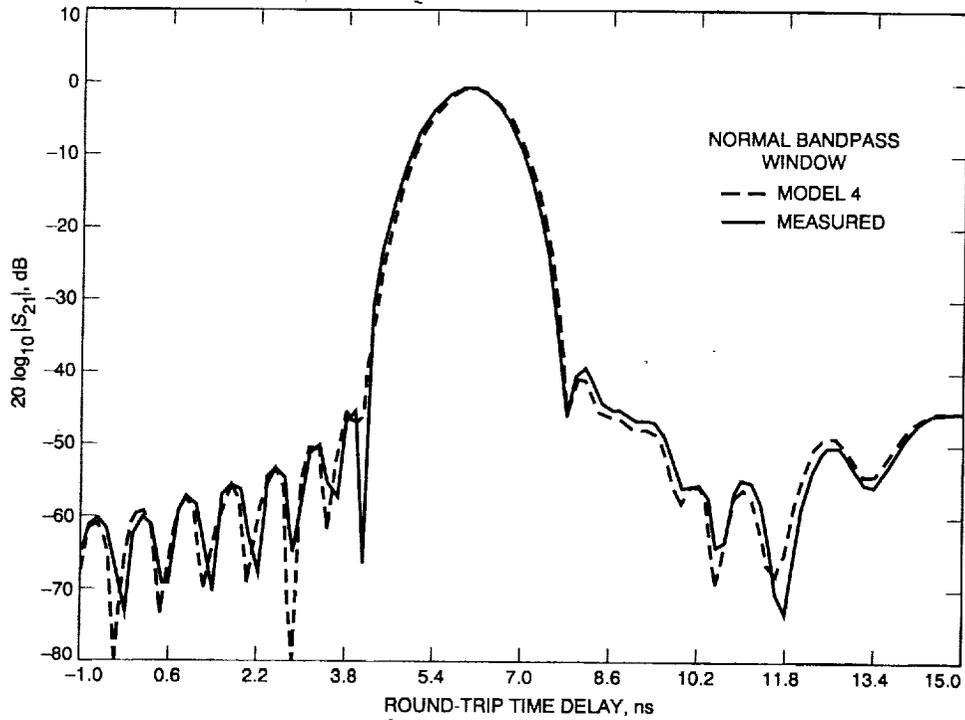


Fig. 14. Comparison of theoretical and measured time domain responses on $|S_{21}|$ data for Model 4.

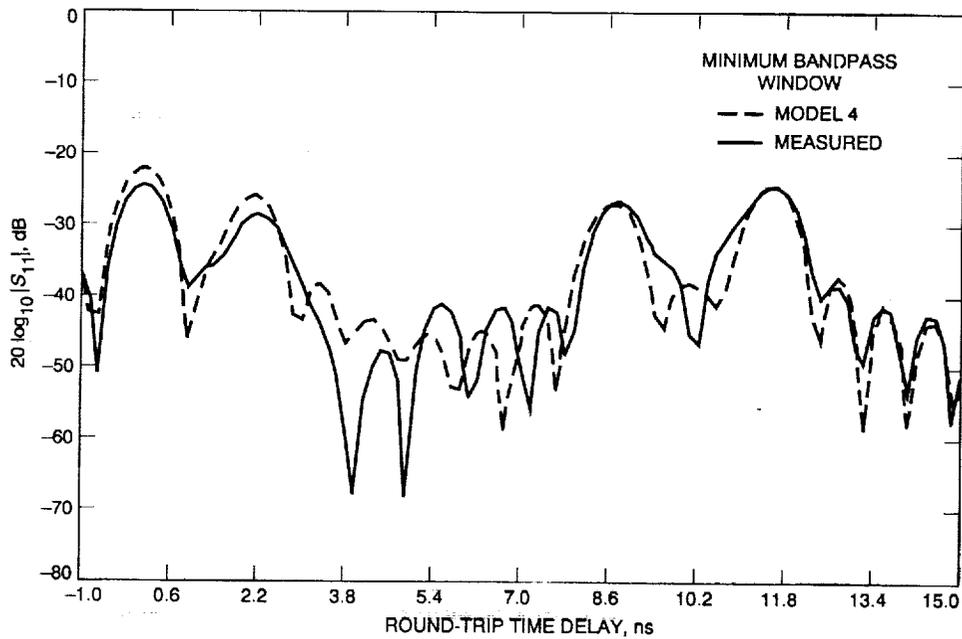


Fig. 15. Comparison of theoretical and measured time domain responses based on $|S_{11}|$ data for Model 4. The theoretical response curve had to be moved by the equivalent of 1.081 in. (0.212 nsec, round-trip) to the left in order to line up the peaks of curves.

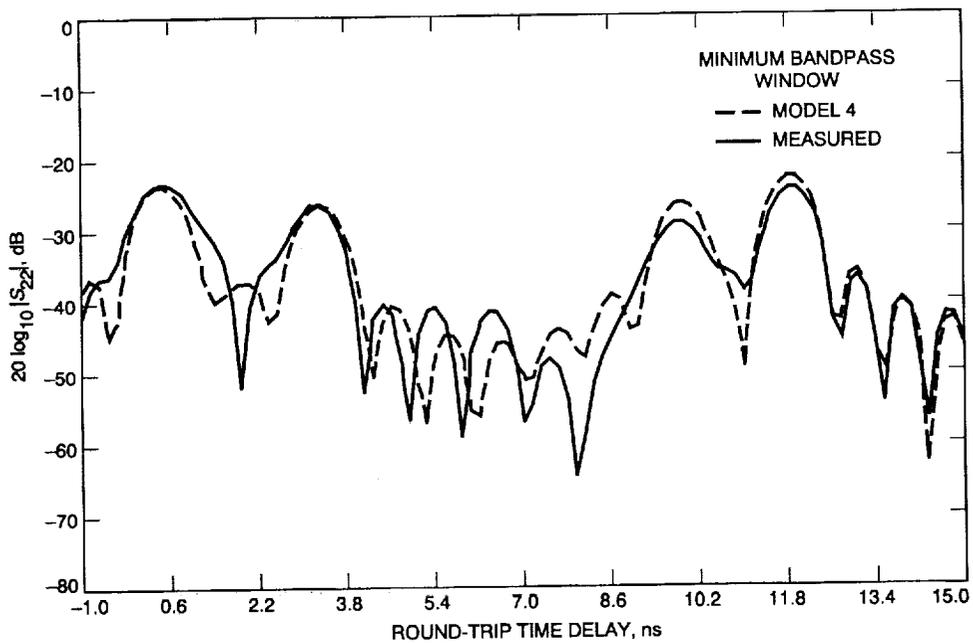


Fig. 16. Comparison of theoretical and measured time domain responses based on $|S_{22}|$ data for Model 4. The theoretical response curve had to be moved by the equivalent of 1.376 in. (0.270 nsec, round-trip) to the right in order to line up the peaks of curves.

Appendix A

S-Parameters of a Basic Network

The basic network used in the modeling work is shown in Fig. A-1. The elements of the network are a capacitive shunt susceptance followed by a length of lossy transmission line. The S -parameters for this network are

$$[S] = \frac{1}{2 + jb} \begin{bmatrix} -jb & 2e^{-\gamma\ell} \\ 2e^{-\gamma\ell} & -jbe^{-2\gamma\ell} \end{bmatrix} \quad (\text{A-1})$$

where b is the normalized shunt susceptance, γ is the complex propagation constant and ℓ is the line length.

For the *Constant Capacitance Model*

$$b = 2\pi f C Z_0 \quad (\text{A-2})$$

where f = frequency in hertz, C is the capacitance in farads, and Z_0 is the transmission line characteristic impedance in ohms.

For the *Constant Shunt Susceptance Model*, b = a constant. The nominal value of b for both models can be obtained from experimental return loss-time domain plots by using the relationship

$$b = \frac{2|S_{11}|}{\sqrt{1 - |S_{11}|^2}} \quad (\text{A-3})$$

$$|S_{11}| = 10^{-(RL_1/20 - A_{dB}/10)} \quad (\text{A-4})$$

where RL_1 is the positive decibel return loss value measured at port 1 and A_{dB} is the positive decibel value for the transmission line loss between the discontinuity and port 1.

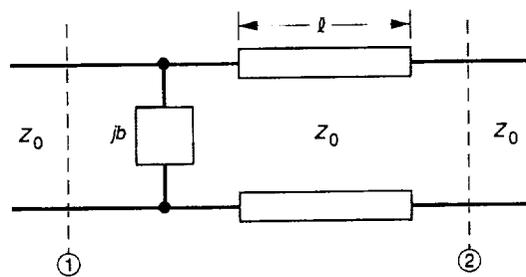


Fig. A-1. Basic network used in cable model.

Appendix B

A Method for Determining Cable-Only Attenuation From Measurements Made on a Cable Having Mismatched Connectors

When a cable has connectors, the method normally used to determine the attenuation of a cable (minus its connectors) is to calculate it from the slope of the measured insertion loss versus frequency curve. When the connectors are mismatched, this method yields an insertion loss that is higher than the cable (minus connector) attenuation. Accurate determination of cable attenuation was critical for the modeling work that was done to obtain the results presented in this article.

From Table B-1, whose equations were derived from those given in [3],

$$|\ell_{11}| \approx \frac{|S_{11}|_{\max} + |S_{11}|_{\min}}{2} \quad (\text{B-1})$$

$$|n_{22}| \approx \frac{|S_{22}|_{\max} + |S_{22}|_{\min}}{2} \quad (\text{B-2})$$

$$|S_{21}|_{\max} = \frac{\sqrt{(1 - |\ell_{11}|^2)(1 - |n_{22}|^2)}}{1 - h} e^{-\alpha l} \quad (\text{B-3})$$

$$|S_{21}|_{\min} = \frac{\sqrt{(1 - |\ell_{11}|^2)(1 - |n_{22}|^2)}}{1 + h} e^{-\alpha l} \quad (\text{B-4})$$

where

$$h = |\ell_{11} n_{22}| e^{-2\alpha l}$$

Then from Eqs. (B-3) and (B-4), the average value of $|S_{21}|$ is derived as

$$\begin{aligned} |S_{21}|_{\text{avg}} &= \frac{|S_{21}|_{\max} + |S_{21}|_{\min}}{2} \\ &= \frac{\sqrt{(1 - |\ell_{11}|^2)(1 - |n_{22}|^2)}}{1 - h^2} e^{-\alpha l} \end{aligned} \quad (\text{B-5})$$

But h^2 is normally less than 0.01 and can be dropped from Eq. (B-5) so that

$$|S_{21}|_{\text{avg}} \approx \sqrt{(1 - |\ell_{11}|^2)(1 - |n_{22}|^2)} e^{-\alpha l} \quad (\text{B-6})$$

then the cable (minus connector) attenuation is

$$\begin{aligned} A_{\text{dB}} &= 20 \log_{10} e^{-\alpha l} \\ &= 20 \log_{10} |S_{21}|_{\text{avg}} \\ &\quad - 10 \log_{10} [(1 - |\ell_{11}|^2)(1 - |n_{22}|^2)] \end{aligned} \quad (\text{B-7})$$

From the pre-environmental cable test data, it was found that near 2 GHz,

$$|\ell_{11}| \approx 0.075$$

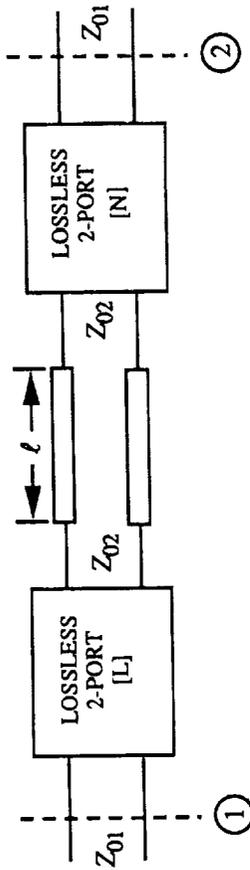
$$|n_{22}| \approx 0.066$$

and the average loss of the cable shown in Fig. 13 between 2.015 and 2.055 GHz is

$$20 \log_{10} |S_{21}|_{\text{avg}} = -0.557 \text{ dB}$$

Substitution of these values into Eq. (B-7) gives A_{dB} equal to -0.514 dB. The contribution of the mismatched connectors from the last term on Eq. (B-7) is about 0.043 dB. For the Galileo cable, the length of the section without connectors is 59.647 in., so the attenuation of the cable (minus connectors) is -0.00861 dB/in. This calculated value compares favorably with the stated manufacturer's data of approximately -0.1 dB/ft or -0.00833 dB/in. However, it was not previously known that the manufacturer's approximate value was close to the true value. Although the corrections seem small, significantly better curve fits between experimental and theoretical data (for Models 3 and 4) were obtained when the more accurate values of the cable attenuation were used at each frequency.

Table B-1. Summary of equations for a lossy uniform line with lossless discontinuities at each end.



Denoting $\arg Z$ as the phase angle of a general-case complex variable Z , let

$$\theta = 2\beta l - \arg \ell_{22} - \arg n_{11}$$

$$\phi = \beta l - \arg \ell_{21} - \arg n_{21}$$

$$h = \ell_{11} n_{22} e^{-2\alpha \ell}$$

$$\tau_0 = \text{delay when } |\ell_{11}| \text{ or } |n_{22}| = 0$$

Parameter	Special cases		
	General case	$\theta = \pm 2(n-1)\pi$	
Overall S-parameters	$S_{11} = e^{j\omega \tau_0} \left[\frac{ \ell_{11} - n_{22} e^{-2\alpha \ell} e^{-j\theta}}{1 - h e^{-j\theta}} \right]$	$\theta = \pm 2(n-1)\pi$	$ S_{11} _{\min} = \left \frac{ \ell_{11} - n_{22} e^{-2\alpha \ell}}{1 - h} \right $
	$S_{12} = S_{21} = \frac{\sqrt{(1 - \ell_{11} ^2)(1 - n_{22} ^2)} e^{-\alpha \ell} e^{-j\phi}}{1 - h e^{-j\theta}}$		$ S_{12} _{\max} = \frac{\sqrt{(1 - \ell_{11} ^2)(1 - n_{22} ^2)} e^{-\alpha \ell}}{1 - h}$
	$S_{22} = e^{j\omega \tau_0} \left[\frac{ n_{22} - \ell_{11} e^{-2\alpha \ell} e^{-j\theta}}{1 - h e^{-j\theta}} \right]$		$ S_{22} _{\min} = \left \frac{ n_{22} - \ell_{11} e^{-2\alpha \ell}}{1 - h} \right $
Group delay ^a	$\tau = \tau_0 \left[\frac{1 - h^2}{1 - 2h \cos \theta + h^2} \right]$		$\tau_{\max} = \tau_0 \left[\frac{1+h}{1-h} \right] = \tau_0 \left[\frac{ S_{21} _{\max}}{ S_{21} _{\min}} \right]$

^a See [3] for derivations.

Appendix C

Equivalent Circuits of a Cable With a Reduced Outer Diameter Section

The equivalent circuits of the section of coaxial transmission line with a reduced outer diameter over length ℓ are shown in Fig. C-1. The characteristic impedance of the nominal and reduced sections Z_{01} and Z_{02} , respectively, are

$$Z_{01} = \frac{60}{\sqrt{\epsilon'_1}} \ell n \left(\frac{D_{O1}}{D_{I1}} \right) \quad (C-1)$$

$$Z_{02} = \frac{60}{\sqrt{\epsilon'_2}} \ell n \left(\frac{D_{O2}}{D_{I2}} \right) \quad (C-2)$$

where D_{O1} and D_{O2} are the diameters of the outer conductors for the nominal and reduced sections, respectively. The symbols D_{I1} and D_{I2} are the diameters of the inner

conductors for the nominal and reduced sections, respectively, and ϵ'_1 and ϵ'_2 are the relative dielectric constants of the media in the nominal and reduced sections, respectively.

Using the equations for the equivalent circuit shown in Fig. C-1(b) given by Beatty [4], the value of $|S_{11}|$ can be calculated. That value is then used in Eq. (A-3) to compute an equivalent shunt susceptance corresponding to the equivalent circuit shown in Fig. C-2(c). If length ℓ is very short (<0.005 wavelength), representing a crimp or deep crease in the outer cable, then the discontinuity should be represented as a shunt capacitive susceptance b of constant value over the frequency range of interest. If the length ℓ is about (0.1–0.25 wavelength), or about equal to the width of a cable clamp, then the equivalent b should be represented as a shunt susceptance with a capacitance of constant value.

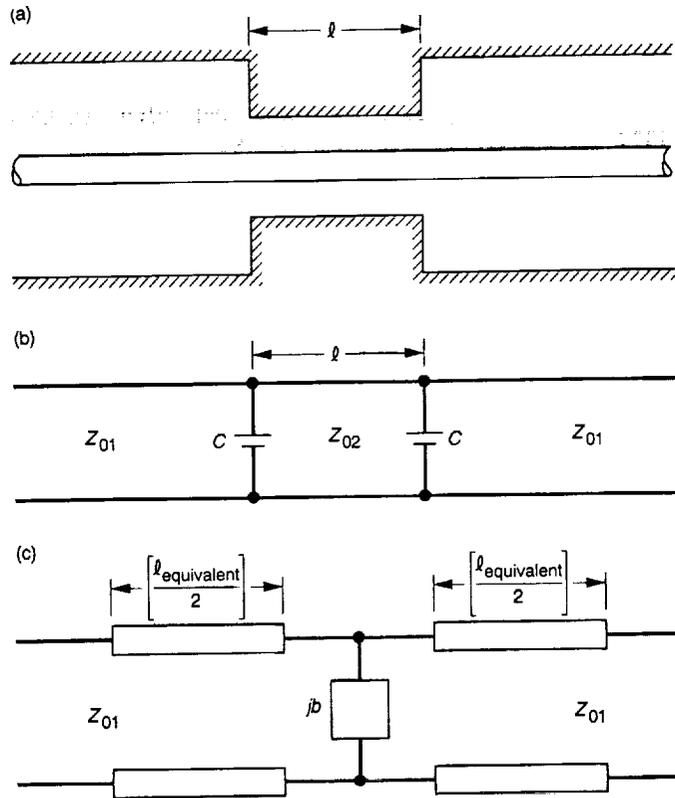


Fig. C-1. Cable with a section of reduced outer diameter: (a) physical representation; (b) equivalent circuit with two capacitive shunt discontinuities separated by a line length of reduced section; and (c) equivalent circuit with a single shunt discontinuity and equivalent line lengths.

References

- [1] C. L. Lawson, "Nonlinear Least-Squares—Plain and Fancy," *Computing and Information Services News*, Jet Propulsion Laboratory, Pasadena, California, vol. 8, no. 4, pp. 6-7, April 1990.
- [2] D. M. Kerns and R. W. Beatty, *Basic Theory of Waveguides and Introductory Microwave Network Analysis*, New York: Pergammon, pp. 42-43, 1967.
- [3] R. W. Beatty and T. Y. Otoshi, "Effect of Discontinuities on the Group Delay of a Microwave Transmission Line," *IEEE Trans. Microwave Theory and Techniques*, vol. MTT-23, pp. 919-922, November 1975.
- [4] R. W. Beatty, "Calculated and Measured S_{11} , and S_{21} , and Group Delay for Simple Types of Coaxial and Rectangular Waveguide 2-Port Standards," *NBS Technical Note 657*, National Bureau of Standards, Boulder, Colorado, p. 18, December 1974.

011-32
128444
N93-19424

DSS-13 Beam Waveguide Antenna Frequency Stability

T. Y. Otoshi and M. M. Franco
Ground Antennas and Facilities Engineering Section

Measurements made on the frequency stability of the DSS-13 34-m-diameter Beam Waveguide (BWG) antenna showed that at 46.5- and 37-deg elevation angles, the BWG antenna stability at 12.2 GHz was between 1.3 and 2.2×10^{-15} for $\tau = 1024$ sec and good weather conditions. These frequency stability values apply to the portion of the antenna that includes the main reflector, subreflector, tripod legs, and the six BWG mirrors. The test results reported in this article are believed to be the first known successful measurements of the stability of the microwave optics portion of a large antenna to a level of 1 or 2 parts in 10^{15} .

I. Introduction

As was pointed out in a previous article [1], attempts have been made in past years to measure the frequency stability of a large antenna using various techniques including (1) a probe on the reflector surface method, (2) a spacecraft Doppler measurement method, (3) a collimation tower method, and (4) very long baseline interferometry (VLBI) methods. Most of these methods have been unable to measure frequency stabilities to better than a few parts in 10^{14} or had other disadvantages discussed in [1].

A new method proposed in 1991 [1] involved the reception of far-field signals from geostationary satellites positioned at various elevation angles. The proposed method had the advantages of being simple and inexpensive to implement and could enable stability data to be obtained in a short time frame. Except for a new fiber-optic subsystem, most of the components and instruments required were already available. Since the antenna now remains the limiting factor that would prevent successful gravitational wave experiments to be performed in the near future, there existed an urgency to obtain antenna stability data that could help establish realistic performance requirements for new DSN antennas.

In addition to the goal of obtaining data in a short time frame, another goal of the new method was to demonstrate that fiber-optic cables could be used to carry microwave frequencies over long distances with negligible degradation to amplitude and phase stability. This article presents data that demonstrates the new method was successfully employed, and that all of the primary goals have been met for good weather conditions. BWG antenna stability data, applicable for inclement weather conditions, are not yet available.

II. Methodology

Figure 1 shows a block diagram of the test configuration. The method involves the use of far-field signals in the 11.7–12.2-GHz region from geostationary satellites, a stable reference antenna, and a phase detector Allan deviation measurement instrument. One of the main advantages of the test method is that it does not require reference signals that are coherent with the station clock. By receiving the far-field signals simultaneously with a reference antenna and the 34-m antenna under test, the phase variations common to both paths tend to cancel at the

output of a mixer contained in the Allan deviation measurement instrument [2].

For the proposed BWG antenna stability measurement method to yield useful and accurate data to the 1×10^{-15} level, it is required that the 12-GHz reference path have a fractional frequency stability of better than 1 or 2×10^{-16} for $\tau = 1000$ sec. The fulfillment of this requirement by the fiber-optic system has been reported in [3].

A Ku-band test package [4] is installed at the pedestal room focal point F3. For this test configuration, the portions of the BWG antenna being tested are the instabilities of the main reflector, subreflector, tripod legs, and six mirrors of the BWG system.

Specific satellites selected for this method are positioned at 47-, 37-, and 12-deg elevation angles. Minor antenna-pointing corrections, of about ± 50 mdeg maximum during a 24-hr period, are required to keep the antenna beam peak pointed at the satellite. Although these antenna pointing changes are small compared with those involved in an actual spacecraft track, phase change measurements could uncover problems associated with antenna pointing due to hardware or software. The method is intended primarily to measure instabilities of the BWG antenna due to outside air temperature and wind conditions.

Even though the measurements are made at Ku-band, the information can be inferred back to mechanically related changes and is therefore useful for determining stability at other frequencies.

III. Test Results

Figure 2 shows the block diagram of the test configuration that was used to measure the stability of the BWG antenna at a 46.5-deg elevation angle. The DSS-13 BWG antenna was pointed at a 46.5-deg elevation angle and a 156.9-deg azimuth for receiving a 12.2-GHz beacon signal from the geostationary satellite GSTAR I owned by the GTE Corporation. The signals from the reference antenna via the fiber-optic system and the 34-m antenna under test are connected to the Allan deviation measurement instrument, which contains a microwave phase detector. Band-pass filters (centered at 12.2 GHz) were used to remove unwanted signals.

Figure 3 shows the Allan deviation plots obtained March 4 and 5, 1992 for two 9-hour-long runs. The first was from 4 p.m. to 1 a.m. PST and the second was from 5 a.m. to 2 p.m. PST. The Allan deviations for $\tau = 1024$

sec were 2.16×10^{-15} for the first run and 1.64×10^{-15} for the second run. These results were about a factor of 3 better than expected.

Also shown in Fig. 3 is the fiber-optic-system-only stability value of 1.64×10^{-16} for $\tau = 1024$ sec. The fiber-optic-only path is the baseline reference that sets the lower limit of stability that can be measured for the 34-m BWG antenna. Descriptions of the fiber-optic-only path and test procedures were described previously in [3].

For interest, Fig. 4 is presented to show the raw phase data corresponding to the Allan deviation result of the above 5 a.m.-2 p.m. run. In addition, the outside air temperature and wind conditions that prevailed during the tests are shown in Figs. 5 and 6, respectively. It can be seen from Fig. 6 that the air temperature varied from about 6 to 15 deg C and the wind was typically less than 30 km/hr. A graphic description of the wind data (Fig. 7) shows that the wind was blowing both into the back and into the front of the main reflector surface at angles of about 16 deg (back side) and 52 deg (front side) off the main z-axis direction.

Figure 8 shows the block diagram of the test configuration used to measure the stability of the BWG antenna at a 37-deg elevation angle. The DSS-13 BWG antenna was pointed at a 37-deg elevation angle and a 132.9-deg azimuth so as to receive a 12.2-GHz beacon signal from the geostationary satellite Satcom K1 owned by the GE American Communications Corporation. This test configuration differs slightly from that shown in Fig. 2 in that one additional 12.198-GHz filter and 20-dB gain amplifier were used. For these tests and all subsequent tests, a new calibration procedure was developed for verification of the test setup. The calibration procedure was to move the subreflector in ± 0.1 -in. offset in the z-axis axial direction from the nominal setting. It is known from previous work [5], that the effective pathlength change was approximately 1.77 times the z-axis offset subreflection position on a 64-m Cassegrain antenna. At 12.2 GHz, this pathlength change, Δ_{pl} , corresponds to a phase change in degrees of

$$\delta_{ph} = \frac{360}{\lambda} (1.77 \Delta_{pl}) \quad (1)$$

where λ = free-space wavelength in centimeters. For $\Delta_{pl} = 0.254$ cm (0.1 in.), then $\delta_{ph} = 65.9$ deg. With the 34-m BWG antenna having shaped main- and subreflector surfaces, the 1.77 factor might be closer to 1.7 so that $\delta_{ph} = 63$ deg. The measured phase changes resulting from moving the subreflector ± 0.1 in. was about ± 60 deg

as shown in Fig. 9. This procedure provided a means of verifying that, for a particular test configuration, phase changes that occurred in the antenna path above F1 were actually observed and correctly measured.

Figure 10 shows the Allan deviation plot obtained when the DSS-13 BWG antenna was pointed at a 37.0-deg elevation angle. The Allan deviations measured for $\tau = 1024$ sec was 1.26×10^{-15} for a 9-hour time period between 4 p.m. to 1 a.m. the next morning.

For interest, the raw phase data for the first run between 4 p.m.-1 a.m. are shown in Fig. 11 along with outside air temperature and wind data in Figs. 12 and 13, respectively. It can be seen from these figures that the air temperature varied between 12 and 7 deg C and the wind was typically less than 24 km/hr. As depicted in Fig. 14, the wind was blowing into the face of the main reflector surface, but at a direction between 29 to 53 deg off the main reflector z-axis.

IV. Future Test Plans

Comparisons of the results showed that for the 4 p.m.-11 a.m. runs, the Allan deviations were 2.16×10^{-15} and

1.26×10^{-15} for $\tau = 1024$ sec at a 46.5- and 37-deg elevation angles, respectively. For future tests, it would be of interest to perform 9-hour runs between 8 p.m. and 5 a.m. when the least amount of air temperature variations generally occur.

For the future, plans are being made to obtain data at the 12-deg elevation as well. It is also of interest to obtain data for the more severe (>32 km/hr) wind conditions that were shown in Figs. 6 and 13.

V. Conclusions

The initial test results presented in this article show that the proposed methodology was successfully employed and the goals were met. Data have been provided to the gravitational wave experimenters in a relatively short time frame (about 14 months) from conception of a new method, followed by the procurement, fabrication, and installation phases, and then the test phase. More data still have to be obtained to determine the stability of the BWG antenna under more severe weather conditions. However, this article has presented the first known stability data obtained on a large microwave antenna to a level of 1 or 2 parts in $\times 10^{15}$ for τ of approximately 1000 sec.

Acknowledgments

D. Bathker of the Ground Antennas and Facilities Engineering Section provided technical and management support on this project, as well as a needed Ku-band feed design for the reference antenna. B. Seidel provided technical information on Ku-band geostationary satellites. Technical discussions with both D. Bathker and B. Seidel concerning the feasibility of the far-field methodology and test strategies were helpful.

The authors are grateful to A. Bhanji, B. Conroy, and R. Perez of the Radio Frequency and Microwave Subsystems Section for the loan of the Allan Deviation Measurement instrument, as well as the technical information on the use of the instrument and data reduction software.

The reference antenna and 12-GHz fiber-optic system were key components that helped to make the measurements successful. J. Garnica, J. Ney, and L. Smith of DSS 13 did most of the work on the fabrication, installation, and checkout of the 10-ft reference antenna. The successful operation of the 12-GHz fiber-optic system was primarily due to G. Lutes and R. Logan of the Communications Systems Research Section.

The support provided by the entire DSS-13 crew during all phases of the testing at DSS 13 is greatly appreciated. C. Goodson and G. Bury provided the DSS-13 management support needed to work out workforce rescheduling problems. Frequent periods of rain and bad weather caused some tests to be cancelled, but due to the excellent station support, the needed test data were ultimately obtained.

References

- [1] T. Y. Otoshi, "A Proposed Far-Field Method for Frequency-Stability Measurements on the DSS 13 Beam-Waveguide Antenna," *TDA Progress Report 42-107*, vol. July-September 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 81-87, November 15, 1991.
- [2] B. L. Conroy and D. Le, "Measurement of Allan Variance and Phase Noise at Fractions of a Millihertz," *Rev. Sci. Instruments*, vol. 61, no. 6, pp. 1720-1723, June 1990.
- [3] T. Y. Otoshi, M. M. Franco, and G. F. Lutes, "Performance of a 12-GHz Fiber-Optic System for Beam-Waveguide Antenna Stability Testing," *TDA Progress Report 42-109*, vol. January-March 1992, Jet Propulsion Laboratory, Pasadena, California, pp. 105-113, May 15, 1992.
- [4] T. Y. Otoshi, S. R. Stewart, and M. M. Franco, "A Portable Ku-Band Front-End Test Package for Beam-Waveguide Antenna Performance Evaluation," *TDA Progress Report 42-107*, vol. July-September 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 73-80, November 15, 1991.
- [5] T. Y. Otoshi and L. E. Young, "An Experimental Investigation of the Changes of VLBI Time Delays Due to Antenna Structural Deformations," *TDA Progress Report 42-68*, vol. October-December 1981, Jet Propulsion Laboratory, Pasadena, California, pp. 8-16, April 15, 1982.

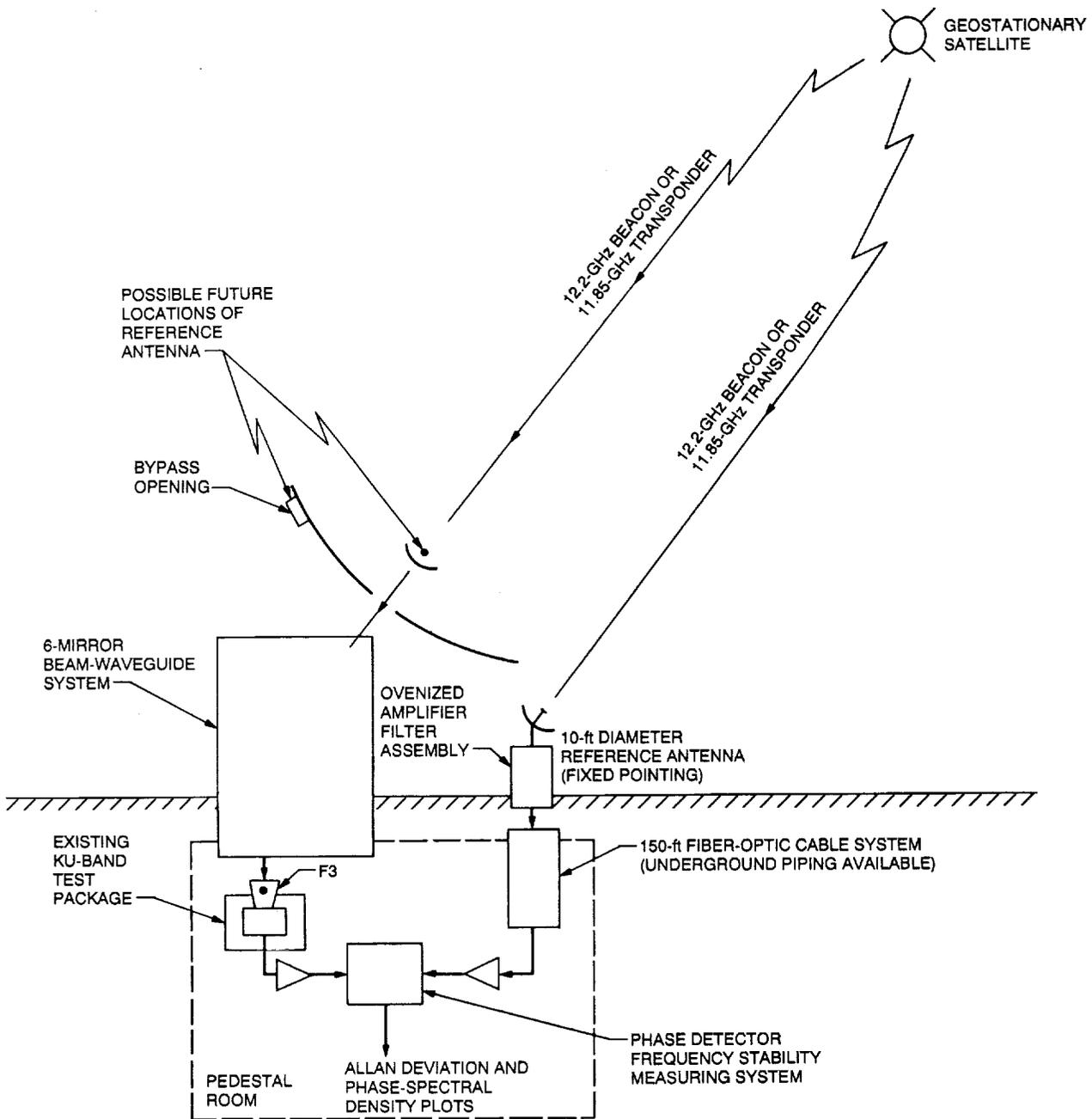


Fig. 1. Test configuration for measurement of the frequency stability of the DSS-13 BWG antenna.

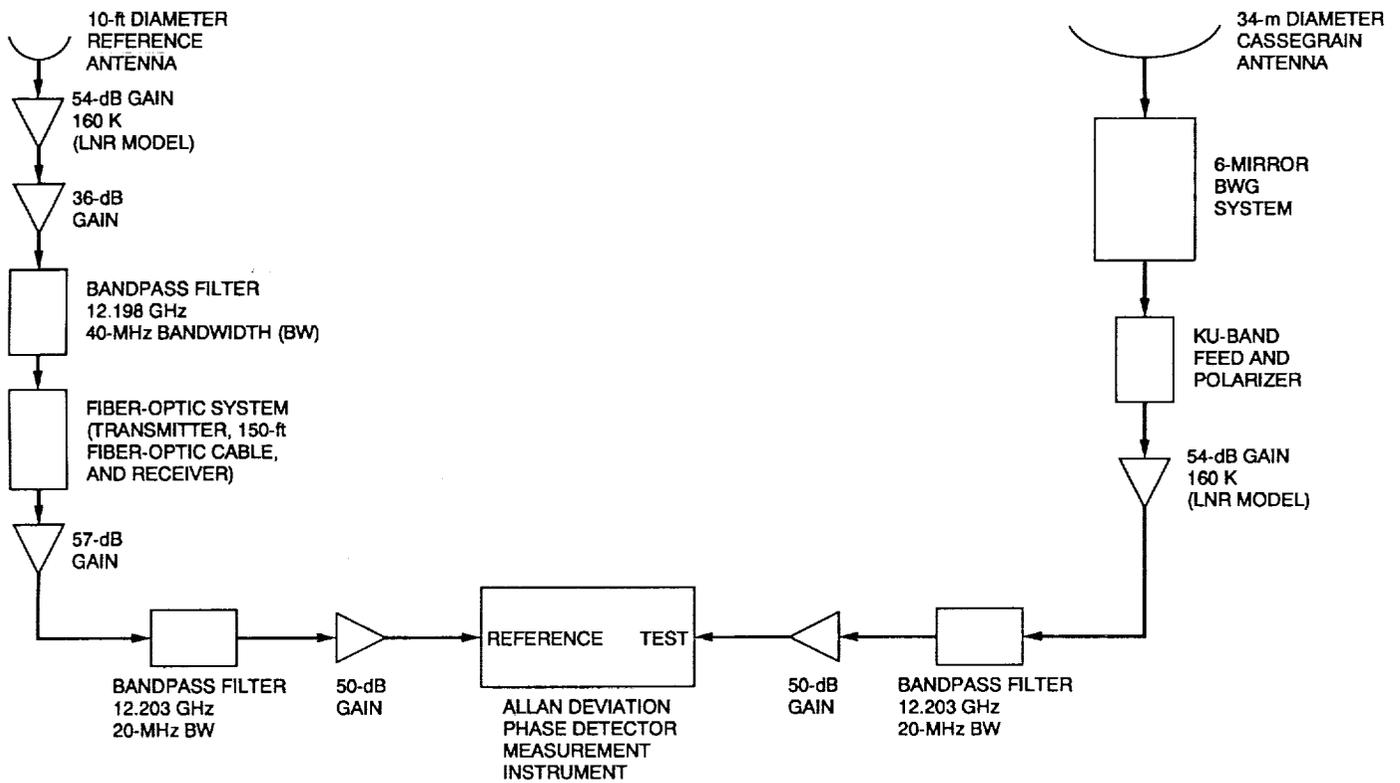


Fig. 2. The test configuration involving the use of the 12.2-GHz beacon signal from the GSTAR I geostationary satellite at a 46.5-deg elevation angle.

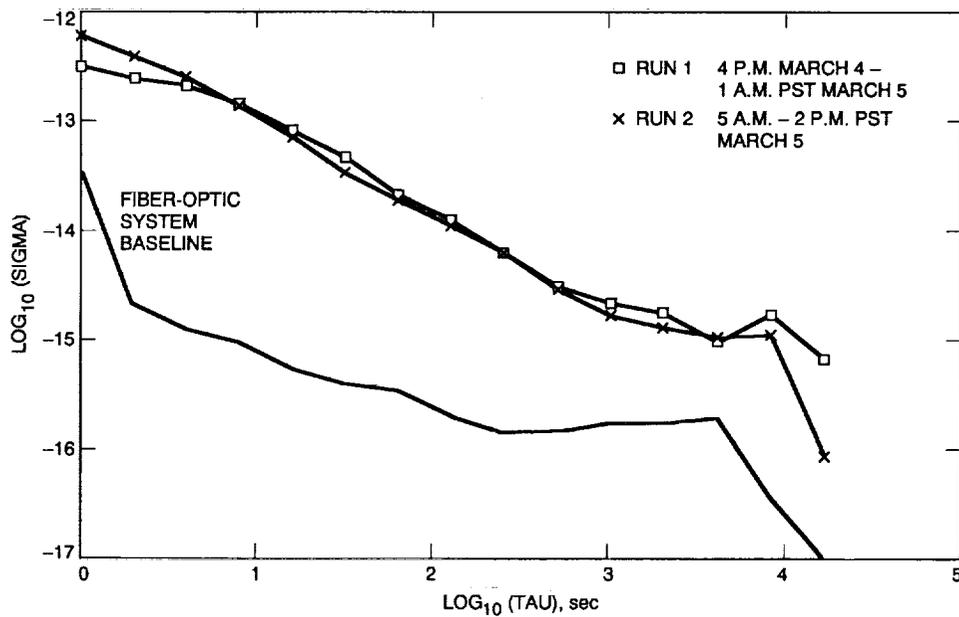


Fig. 3. Allan deviation plots of the DSS-13 BWG antenna stability on March 4 and 5, 1992. The DSS-13 BWG antenna was pointed at a 46.5-deg elevation angle and a 156.9-deg azimuth, receiving a 12.2-GHz beacon signal from GSTAR I.

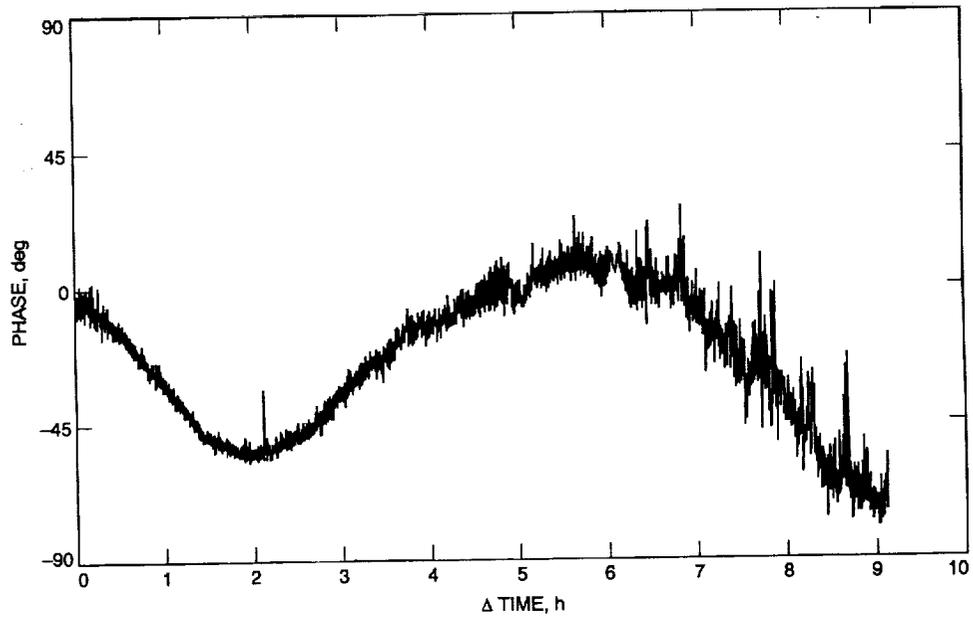


Fig. 4. Phase change plot corresponding to the 5 a.m.-2 p.m. run of Fig. 3.

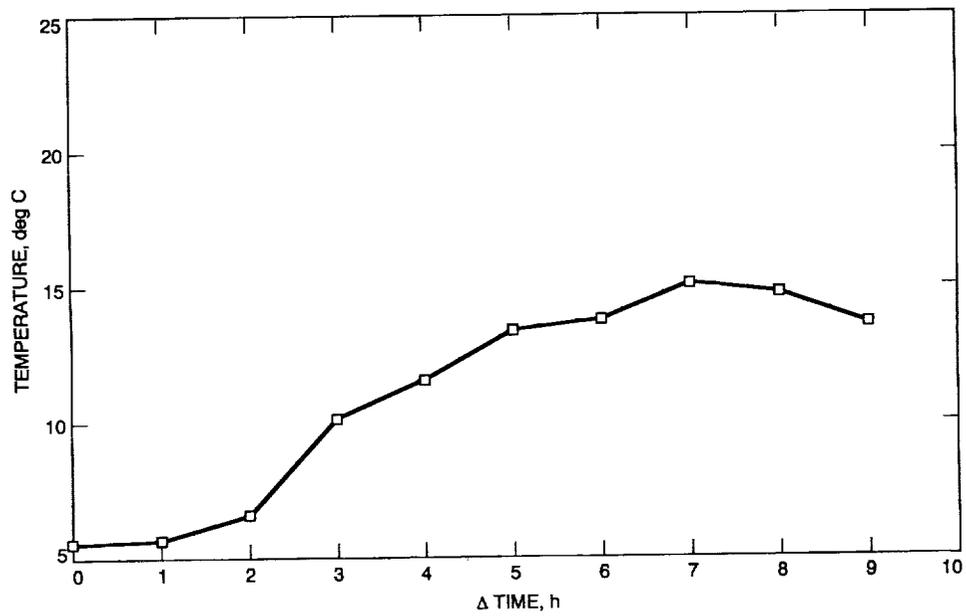


Fig. 5. Outdoor ambient temperature change during the 5 a.m.-2 p.m. run of Figs. 3 and 4.

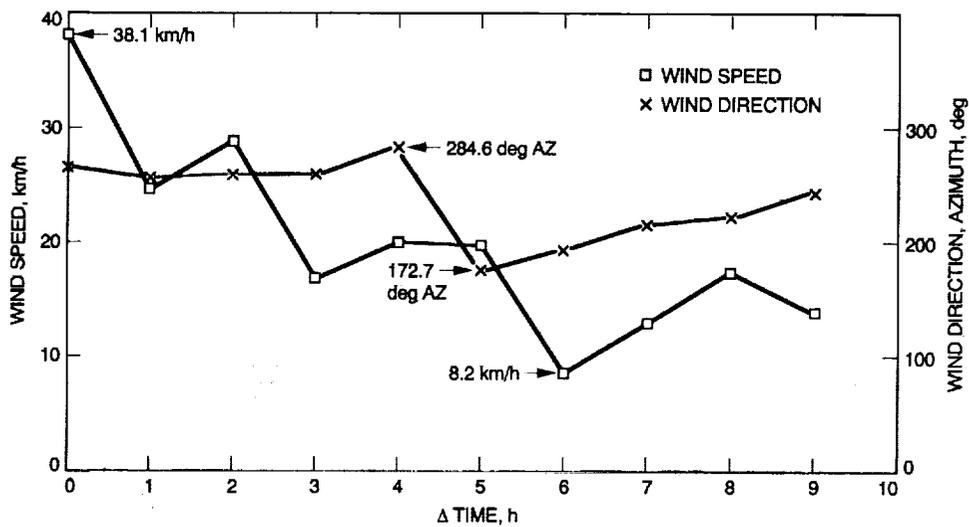


Fig. 6. Outdoor wind data during the 5 a.m.–2 p.m. run of Figs. 3 and 4.

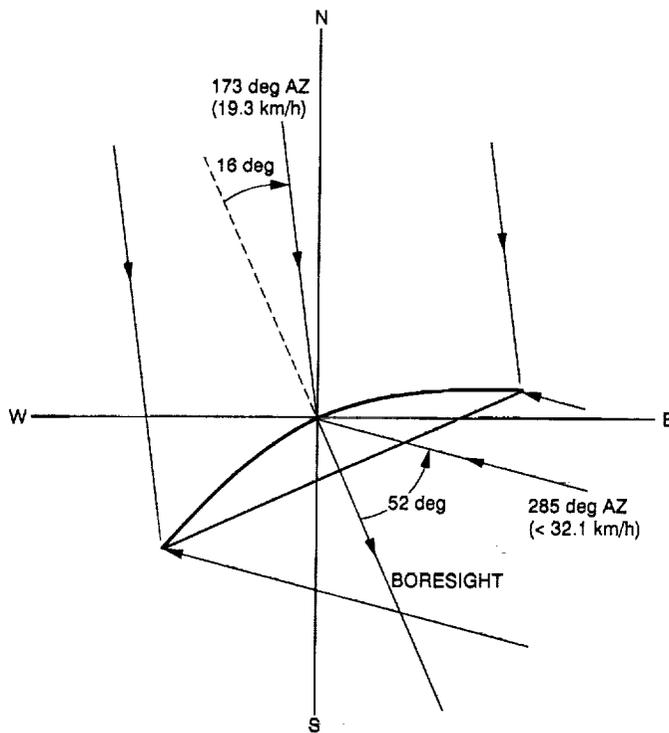


Fig. 7. Wind direction relative to the BWG antenna main axis during tests with the GSTAR I geostationary satellite.

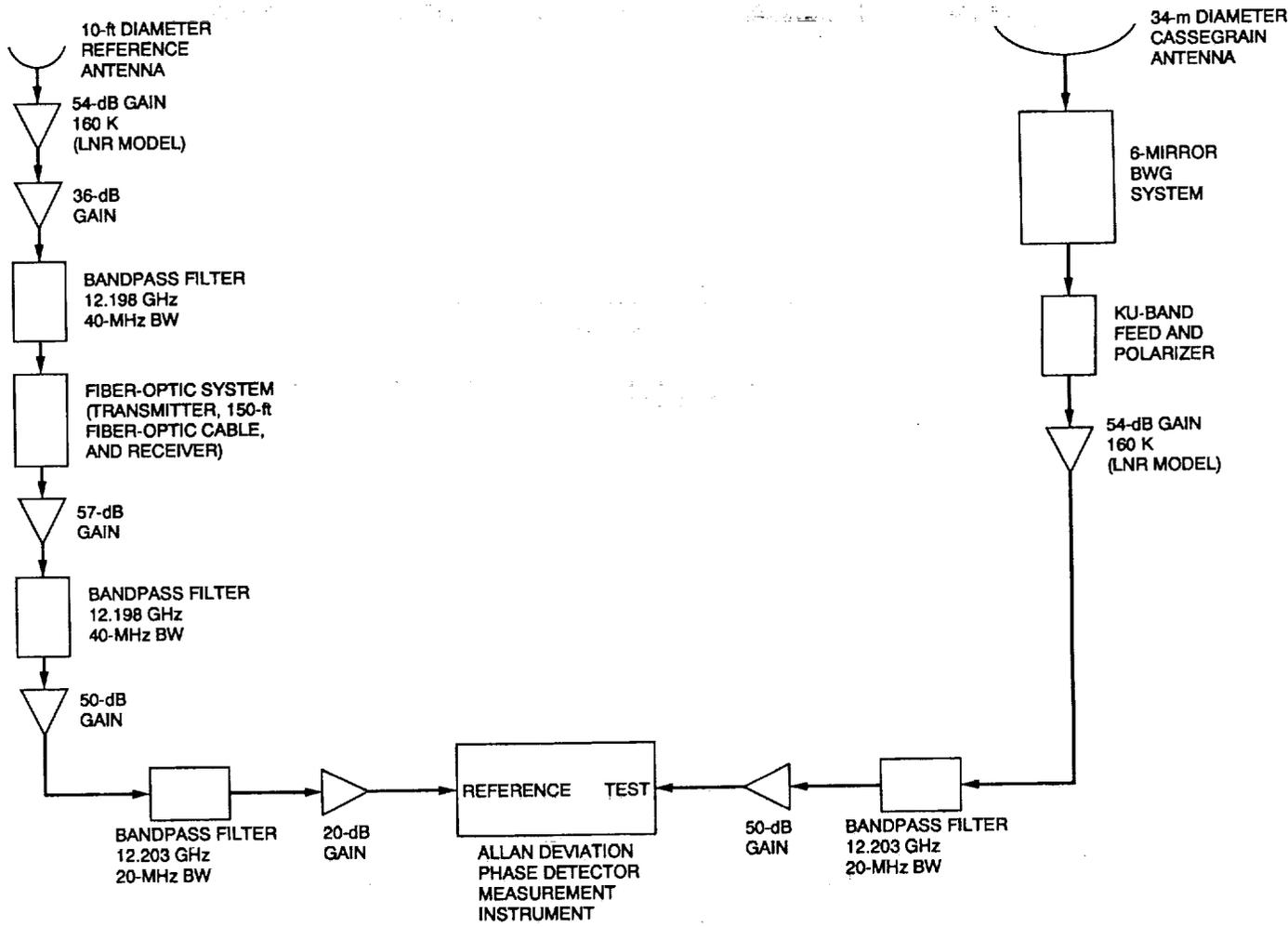


Fig. 8. The test configuration involving the use of the 12.2-GHz beacon signal from the Satcom K1 geostationary satellite at a 37-deg elevation angle.

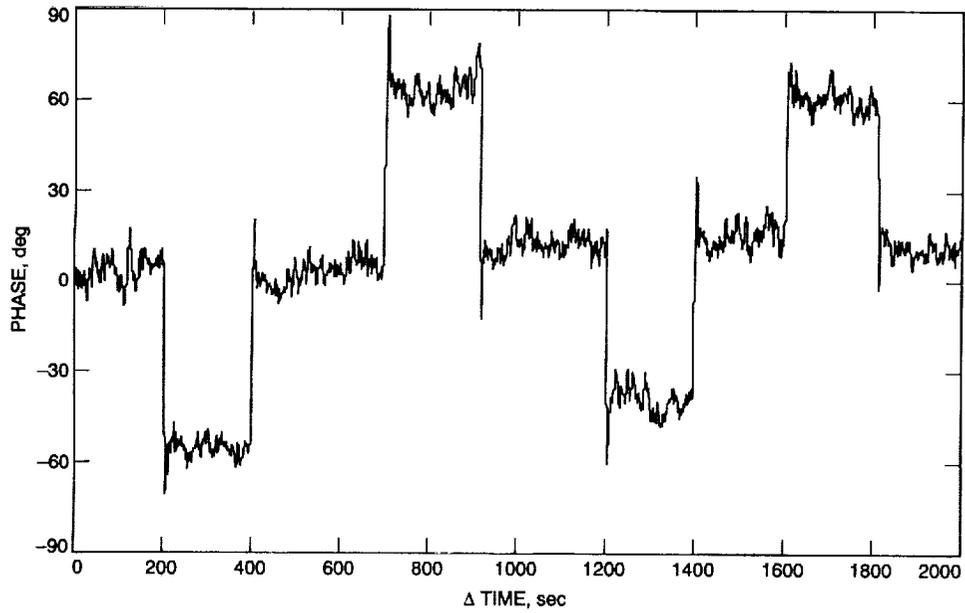


Fig. 9. Subreflector phase change plot as part of the precalibration procedure to verify test setup on 1992 day of year 66. The DSS-13 BWG antenna was pointed at a 37-deg elevation angle and a 132.9-deg azimuth, receiving a 12.2-GHz beacon signal from Satcom K1.

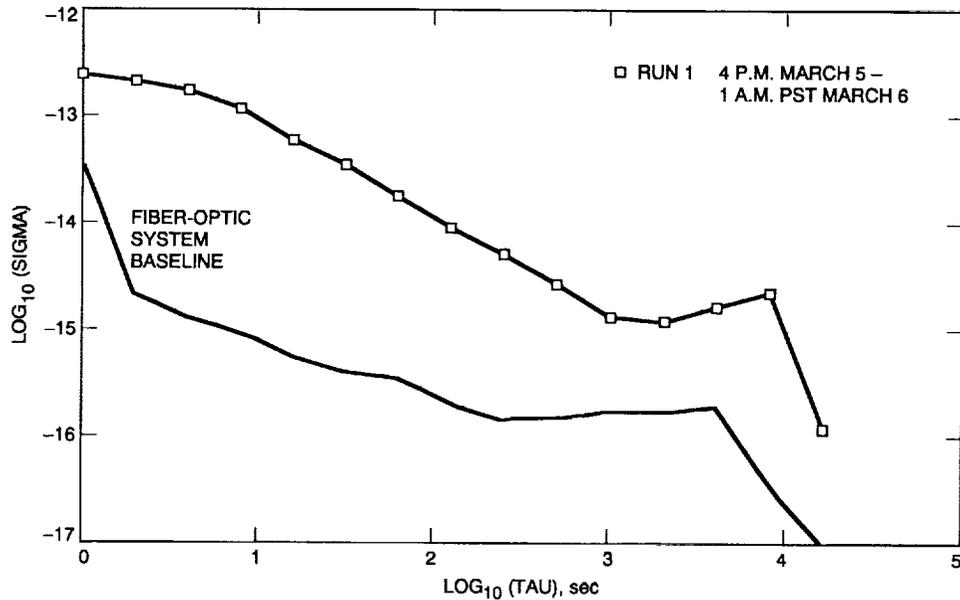


Fig. 10. Allan deviation plot of the DSS-13 BWG antenna stability on March 5 and 6, 1992. The DSS-13 BWG antenna was pointed at a 37-deg elevation angle and a 132.9-deg azimuth, receiving a 12.2-GHz beacon signal from Satcom K1.

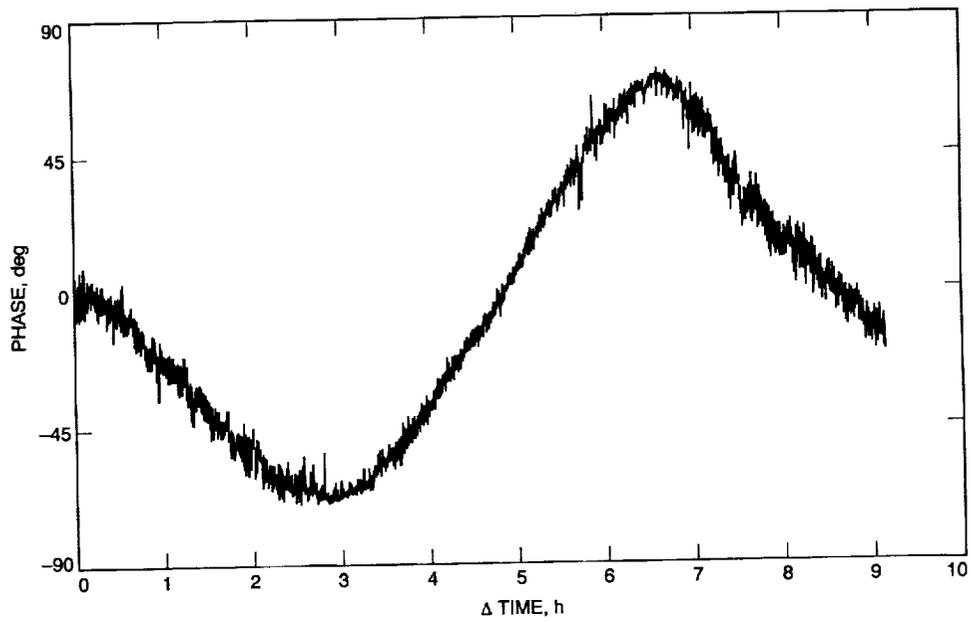


Fig. 11. Phase change plot corresponding to the 4 p.m.-1 a.m. run of Fig. 10.

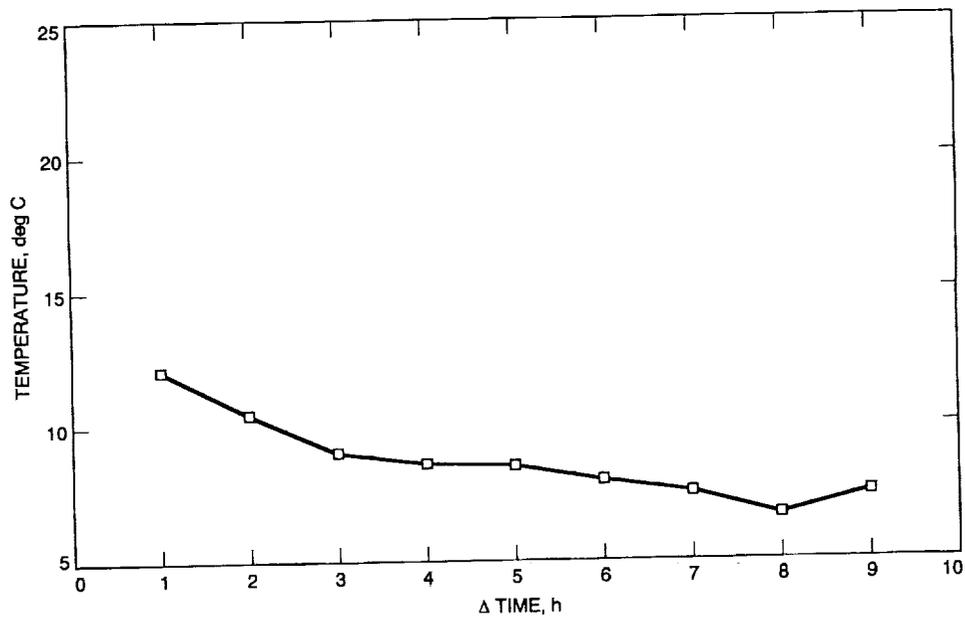


Fig. 12. Outdoor ambient temperature change during the 4 p.m.-1 a.m. run of Figs. 10 and 11.

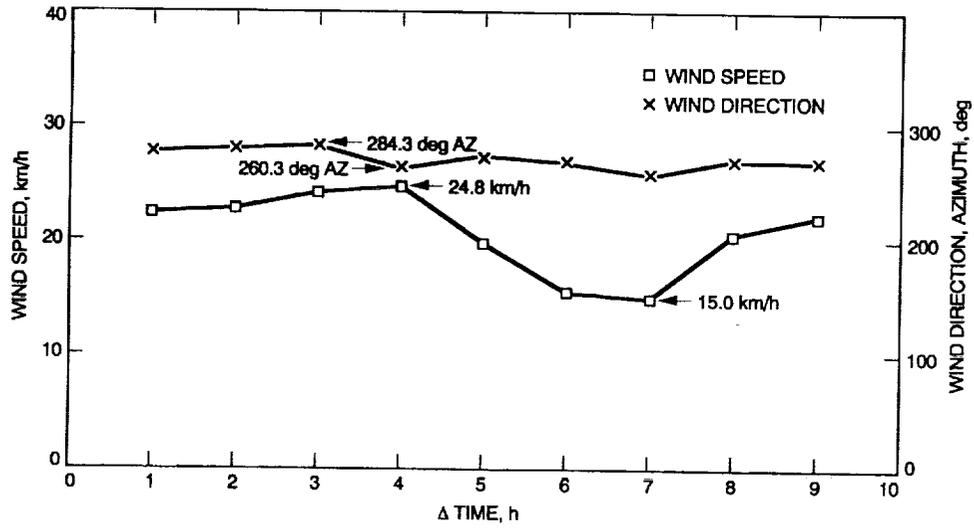


Fig. 13. Outdoor wind data during the 4 p.m.–1 a.m. run of Figs. 10 and 11.

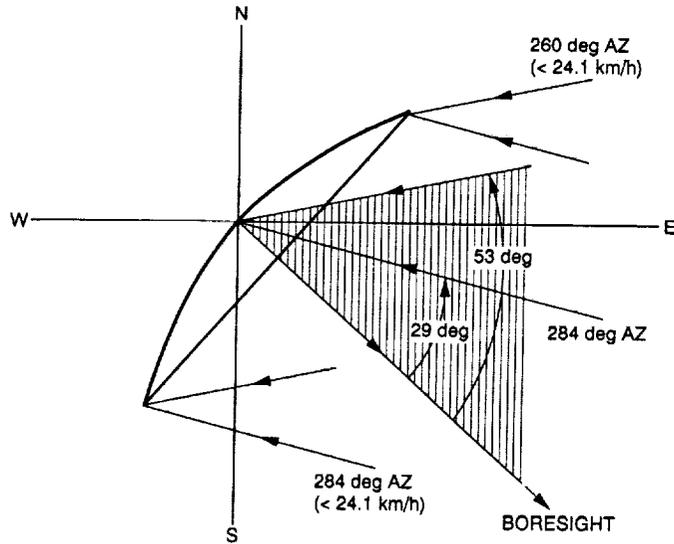


Fig. 14. Wind direction relative to the BWG antenna main axis during tests with the Satcom K1 geostationary satellite.

N 93 - 19425
 513 61
 128445
 p-16

Locally Adaptive Vector Quantization: Data Compression With Feature Preservation

K.-M. Cheung

Communications Systems Research Section

M. Sayano

California Institute of Technology

This article presents a study of a locally adaptive vector quantization (LAVQ) algorithm for data compression. This algorithm provides high-speed one-pass compression and is fully adaptable to any data source and does not require a priori knowledge of the source statistics. Therefore, LAVQ is a universal data compression algorithm. The basic algorithm and several modifications to improve performance are discussed. These modifications are nonlinear quantization, coarse quantization of the codebook, and lossless compression of the output. Performance of LAVQ on various images using irreversible (lossy) coding is comparable to that of the Linde-Buzo-Gray algorithm, but LAVQ has a much higher speed; thus this algorithm has potential for real-time video compression. Unlike most other image compression algorithms, LAVQ preserves fine detail in images. LAVQ's performance as a lossless data compression algorithm is comparable to that of Lempel-Ziv-based algorithms, but LAVQ uses far less memory during the coding process.

I. Introduction

Data compression is the art of packing data, the process of transforming a body of data to a smaller representation from which the original or an approximation to the original can be computed at a later time. Most data sources contain redundancies such as nonuniform symbol distribution, pattern repetition, and positional redundancy. A data compression algorithm encodes the data to reduce these redundancies.

Data compression has not been a standard feature in most communication/storage systems for the following

reasons: Compression increases the software and/or hardware cost; compression/decompression is difficult to incorporate into high data rate (greater than 10 Mb/sec) systems; most compression algorithms are not flexible enough to process different types of data; the unpredictability of compressed data file size presents space allocation problems. These obstacles are less significant today due to recent advances in algorithm development, high-speed very large-scale integrated circuit (VLSI) technology, and packet switching communications. Data compression is now a feasible option for those communication or storage systems for which communication bandwidth and/or storage capacity are at a premium. If present

trends continue, the volume of speech and image data in the near future will become prohibitively large for many communication links or storage devices.

A number of applications require data compression for efficient data storage. To facilitate fast processing, maintain accurate records, and recall old records quickly, a number of businesses are optically scanning their documents and saving them on magnetic media. This takes up large amounts of memory; efficient storage requires data compression. Documents are stored for archival purposes, so a short delay in retrieval (decompression) is not detrimental. A similar application exists in law enforcement, security, and intelligence agencies, where facial, fingerprint, and other images are kept on file for fast retrieval, analysis, and matching.

Data compression is required also in limited bandwidth communications; the two best examples of this are video telephony and high-definition television (HDTV). Video telephony requires transmission of images over a small bandwidth; this can be as small as 4 kHz in the standard voice communication channel. Sending an image, even a small one, without data compression is not feasible. For HDTV, data compression is also needed if signals are to be sent digitally over a standard television channel: HDTV signals require roughly four times the bandwidth of standard TV signals. High-speed algorithms which compress the image without substantially degrading image quality are requirements for both video telephony and HDTV. As more and more information must be transmitted over the same size bandwidth, data compression becomes imperative to maintain transmission rate and data fidelity.

Vector quantization (VQ) is an efficient data compression technique for speech and images. VQ maps a sequence of continuous or discrete vectors into a digital sequence suitable for transmission over a digital channel or storage in a digital medium. The goal is to reduce the volume of data while preserving required fidelity levels. In [1], a well-designed VQ scheme was shown capable of providing high compression ratio with good reconstructed quality.

Unlike scalar quantization where the actual coding of continuous or discrete samples into discrete quantities is done on single samples, the input data of a VQ encoder are multidimensional blocks of data (input vectors). An important technique in VQ is the training of codebooks prior to transmission [1]. Extensive preprocessing is performed on sample source data to construct the codebook to be used in the compression session. The encoder and decoder must first agree on the same codebook before data transmission. The closeness between an input vector and

a codeword in the codebook is measured by an appropriate distortion function. During the compression session, distortions between an input vector and codewords in the codebook are evaluated; the codeword closest to the input vector is chosen as the quantization vector to represent the input vector. The index of this chosen codeword is then transmitted through the channel. Compression is achieved since fewer bits are used to represent the codeword index than the quantized input data. The decoder receives the codeword index and reconstructs the transmitted data using the preselected codebook.

Traditional VQ schemes have a few inherent disadvantages. (1) The generation of a good codebook requires a priori knowledge of the source data, which in practice are not often easily available. (2) Traditional VQ schemes are static schemes. They assume that the statistical properties of the source data remain the same for all compression sessions and the codebooks are optimized based on this assumption. Real-world data tend to have varying characteristics, and static algorithms may not be efficient enough to process diverse sources. (3) Both codebook generation and codeword search for input vectors involve computing the distortion between input vectors and codewords in the codebook; these are usually computationally intensive processes, especially when the codebook size is large.

A new class of data compression algorithms, locally adaptive vector quantization (LAVQ) algorithms, was suggested in [2] and [3]. Unlike traditional VQ algorithms, LAVQ algorithms do not require a priori knowledge of the source, nor do they require the tedious process of codebook generation, as the codebook is generated on the fly during encoding, and the decoder mimics the operations of the encoder to maintain an identical codebook at all times. This algorithm does not require a full codebook search: The codebook is updated after each use to maintain the most recently used codewords at the front of the codebook; a codeword which is within the error allowance is typically found in the top one-fifth to one-tenth of the book through a sequential search. This algorithm dynamically adapts to the local features of the source and is particularly good in compressing sources with varying characteristics.

In this article the basic algorithm suggested in [2] and [3] will be described. Subsequent sections will be devoted to the analogy of LAVQ to a vector differential pulse code modulation (DPCM) algorithm, improvements to the basic LAVQ algorithm, and application to lossless data compression. The algorithm has been fully implemented in software; a brief description of this implementation is included. Experimental results on both lossy image coding and lossless data compression are also presented.

II. Basic Algorithm

The basic LAVQ algorithm provides a simple yet effective one-pass data compression strategy (refer to Fig. 1). The encoder has a codebook containing codewords (vectors) where the index of the codeword corresponds to its position in the codebook. A block is taken from the image and compared to the stored codewords; if there exists a codeword sufficiently close to the image block (within the error allowance), the index itself is sent, and that codeword is moved to the top of the codebook. If no such codeword exists, a special index is sent. This index is followed by the block itself. This block becomes a new codeword and is placed at the top of the codebook. All other codewords are pushed down, and if the number of codewords exceeds the maximum allowed, the last codeword is lost. Initially, the codebook may be empty or full from the previous image encoded.

On the decoder side, the decoder expects an index. If this index is the special one denoting that a new block was sent, the decoder expects a block to be received immediately following the special index; this block is placed at the top of the codebook and all other codewords are pushed down. If the codebook is already full, the last codeword is discarded. This new block is also placed into the image being built by the decoder. If the index is not one designating a new block, then the codeword corresponding to the index is put into the image being built, and that codeword is moved to the top of the codebook. Thus, if the encoder and decoder start with the same codebook, they will have the same codebook at each step, and the image will be successfully sent [2-4].

The LAVQ strategy maintains the most recently used vectors in the codebook in the order of last usage; this allows the algorithm to efficiently code any image on the fly without codebook training: The algorithm needs only one pass of the image to code it entirely. In serial implementation, LAVQ has time complexity $O(nm)$ and storage complexity $O(m)$, where n is the number of pixels in the image and m is the number of codewords in the codebook. Most of the time spent on encoding is taken by finding the closest codeword in the codebook and determining if that match is close enough. Rearranging the codebook and sending the required index, and possibly a new block, can be done quickly in serial implementation using lookup tables, linked lists, and other software techniques. To minimize the amount of time spent on codebook searching and to improve performance, a partial search of the codebook can be used: Instead of searching serially for the best match, the encoder can stop searching at the first instance of a close-enough match, with the criterion dictated by the error allowance given.

The basic LAVQ algorithm, however, has poor performance compared to traditional VQ strategies such as the Linde-Buzo-Gray (LBG). Several adjustments can be made to improve the algorithm without degrading the advantage of one-pass high-speed implementation. Two approaches can be used to improve rate. First, the statistics of the coded indices can be skewed toward small values (recent indices) by (1) a partial codebook search as outlined above, (2) using tall and narrow blocks ($N \times 1$ pixels) to make each block more similar to the blocks immediately previous to it in a raster scan, and (3) coding only the differential value of each block by removing the mean value and reinserting it later. Second, the number of bits used to send new codewords can be reduced by (1) reducing the number of bits used to describe each new pixel and (2) using nonlinear quantization of these new codeword values. These two approaches, combined with lossless adaptive arithmetic coding of the output, improve the performance of LAVQ to be comparable to that of LBG. These improvements will be discussed in detail in a subsequent section.

III. LAVQ as a Vector Analogy of DPCM

Differential pulse code modulation (DPCM) data compression algorithms are efficient and have low complexity. They are particularly effective in encoding gray-scale images, which are dominantly characterized by an autoregressive (AR) stochastic model or an autoregressive moving average (ARMA) model. DPCM operates on individual samples $x(n)$ and encodes the quantized difference $e(n)$ between a predicted value $\hat{x}(n)$ and $x(n)$. The prediction is based on the pixels neighboring $x(n)$. The error $e(n)$ tends to be small rather than large, and compression is achieved by assigning fewer bits to smaller values of $e(n)$ and more bits to larger values of $e(n)$.

From an information theoretic point of view, given a data source, it is always advantageous to encode vector quantities rather than scalar quantities. A vector extension of DPCM coding involves encoding a vector of differences $E = [e_1(n), e_2(n), \dots, e_N(n)]$, where N is the vector size. However, the number of combinations of E increases exponentially with N ; especially for large N , this quickly becomes too large to be practically feasible for efficient encoding. Another drawback of DPCM is its inflexibility: DPCM performs well when coding sources are characterized by the AR or ARMA models (e.g., speech and image data); however, for data sources dominantly characterized by pattern repetitions (e.g., data base records and engineering data), DPCM performs poorly.

LAVQ is analogous to DPCM in a sense: Both DPCM and LAVQ can be viewed as consisting of a preprocessing stage and a compression stage. DPCM preprocesses a sample by taking the difference between the sample value and its corresponding predicted values and sends the quantized difference to be entropy coded. LAVQ, on the other hand, preprocesses a vector of samples by matching it to the codeword vectors in a dynamic codebook followed by a move-to-front codebook update, and sends either a codebook index or an uncoded vector to the decoder. The major difference is that DPCM encodes the scalar difference, whereas LAVQ encodes the vector recency, which can be considered as a different measure of difference (vector difference). Thus, LAVQ uses the locally adaptive move-to-front preprocessing unit to convert the hidden statistical and correlational redundancy to a scalar statistical redundancy. This allows representation of vectors as scalars and approaches vector entropy, which, in the information theoretic sense, is smaller than the corresponding scalar entropy. This was proven in [2] and is verified in the results given here.

IV. Improvements to LAVQ for Image Compression

A. Index Coding

1. Difference Coding. The most significant visual artifact of LAVQ is the sawtooth or staircase effect. This occurs from using discrete vectors to represent a set of vectors within a threshold set by the error allowance; thus, differences between codewords can be large enough to be noticeable. In particular, if the pixels are slowly varying across the image, as is the case with most images, the encoded blocks do not track this variation closely. That is, adjacent blocks are coded with the same block while the amount of error is within the allowance; when that allowance is exceeded, suddenly a different block is used. This difference can be noticeable and, since most images have regions of slow variations or of constant color or intensity, quite common.

This assumption of images having regions of slowly varying or constant pixels implies that adjacent blocks have similar mean values. Thus, the mean can be removed, and only the differential values of the image can be coded. Each vector now represents the difference between the actual pixel value and a reference value equal to the distorted mean of the previously coded block. This distorted mean is a valid choice of reference because both the encoder and decoder can compute it exactly from the previous block's distorted pixel values. The small difference vectors are

more likely to be well approximated by the recently occurring difference vectors maintained in the current codebook. Thus, slowly varying or constant regions can be more accurately coded without extensive use of new codewords. The mean must be updated after each block is processed at both the encoder and the decoder. This ensures that both have the same mean to remove and to reinsert into the image at each step.

2. One-Pass Index Compression. Because of the move-to-front codebook rearrangement strategy and because most images have similar adjacent blocks, the smallest indices are most likely to be used more often. Therefore, they can be coded using a lossless compression code to obtain better performance. However, to maintain one-pass compression, the lossless code must also be one-pass. Furthermore, the statistics of the coded indices are unknown a priori, and no assumption can be made regarding them.

The code which yields the most promising results with minimal increase in computational complexity is the adaptive arithmetic code. This algorithm is the static arithmetic code implemented with probabilities of each symbol updated after each use. By starting with the same initial distribution at the encoder and decoder, lossless coding can be obtained. The arithmetic code approaches global symbol entropy closely; in some cases, it does even better: The average entropy of the adaptive arithmetic code is average local entropy based on symbol statistics from the start to the symbol being coded. Global entropy is derived from the statistics of each symbol based on the entire sequence; local entropy is derived from the statistics of each symbol based on part of the sequence in the neighborhood of the symbol being coded. The global entropy is always greater than or equal to the global average of all the local entropies. Thus, for sources which have localized characteristics which vary throughout the sequence, using localized adaptive coding methods is more advantageous than using global, nonadaptive coding methods.

B. Codeword Data Coding

1. Bit Stripping. Bit stripping of new codeword values can be used independently from or in conjunction with difference coding to obtain higher compression rates. The least-significant bits of either a block of pixels or a vector of differences tend to be uniformly random; therefore, stripping these bits before sending data to the decoder and reinserting a mean value at the decoder decreases the number of bits sent when a new codeword is generated for the codebook, with a small increase in distortion. The amount of additional error incurred is not readily noticeable in most images.

2. Nonlinear Quantization. When difference coding is used with bit stripping, only a few values typically occur, and these may be represented by a relatively small number of quantization levels. However, there may be sharp edges in the image which will have large differences in value between adjacent blocks. This can cause large errors at edges if linearly quantized difference values are too small to keep up with the rapid change in pixel values. To minimize this error, nonlinear quantization can be used. The choice of quantization step sizes is not clearly defined, as the image statistics are unknown to the encoder and cannot be easily transmitted to the decoder. Therefore, design of the quantizer cannot be made adaptive, and an arbitrary choice must be made beforehand. Quantizer design has some criteria, however. Initial step sizes (near zero) should be small, and subsequent sizes should increase. Based on this assumption, a fixed logarithmic quantizer is used.

C. Interpolation and Smoothing

As mentioned earlier, since LAVQ uses discrete vectors to represent a set of vectors within a threshold set by the error allowance, there is noticeable blockiness in the output. Difference coding helps remove this effect, but it is insufficient: Since zeroth order estimation of the next block is used—the mean is assumed to be identical across adjacent blocks—regions with gradual pixel value changes can cause a sawtooth or staircase visual artifact at high compression rates. In addition, difference coding does not remove the block boundaries visible in the vertical direction. In cases where difference coding is used and where it is not, horizontal interpolation can remove the horizontal artifacts, and vertical smoothing can remove the vertical artifacts; if done carefully, both techniques do not excessively destroy detail or cause a blurry appearance.

Horizontal interpolation essentially interpolates across those blocks represented by different codewords. Each block is classified as either a repeat of the previous block or not; the first occurrence of a block that is not a repeat of the previous block is recorded; the rest are initially left blank. These blank blocks are filled with a pixelwise linear interpolation of the edges of the two closest nonblank blocks. However, this can smear edges which occur in the image; therefore, a threshold is used: If the difference between two differing blocks is larger than this threshold, no interpolation is done. This threshold must be adjusted externally.

Vertical smoothing averages the pixels on the vertical block borders. If both blocks on the border are not new codewords, that is, both already exist in the codebook, then the two border pixels are averaged; this average value

is substituted for the border pixels. If both blocks are new codewords, nothing is done. If only one is new and the other is an existing codeword, then the new block is unaltered; the existing block's boundary pixel is substituted with the average of the boundary pixel of the new block and the two boundary pixels (the boundary and the pixel vertically adjacent to it) of the existing block. In this way, new codewords, which usually describe detailed areas, are unaltered, while existing codewords, which describe areas of low detail, are smoothed.

V. Lossless Data Compression

As mentioned earlier, the LAVQ algorithm is also suitable for lossless compression of database records and other data dominantly characterized by pattern repetitions. Examples of these data sources include textual data, accounting and payroll database data, telemetry data, and engineering data. Lossless compression using LAVQ can be achieved easily by setting the error allowance to zero. In the lossless compression mode, the basic LAVQ encoder becomes a locally adaptive move-to-front (MTF) algorithm.

Lossless LAVQ works best on data with fixed record sizes. Data represented in fixed-size packets and with patterns confined to the packets or fixed subfields within them are the best candidates for LAVQ. However, arbitrary fixed-length blocking of data can also be used on other sources without defined block sizes without significant detriment. Universal data compressors such as Lempel-Ziv (LZ)-based algorithms assume no structure of the source except for intersymbol correlation; however, many of these algorithms require large amounts of memory during coding, use complex tree data structures such as the "Patricia tree" to maintain their dictionaries, and have sophisticated pruning techniques to update the code tree. LAVQ requires relatively less memory and less complex data structures; it uses a simple codebook updating algorithm (MTF in this case) with a much smaller data buffer. Several high-performance LZ-based algorithms use back-end entropy coders to further improve performance; LAVQ can be likewise equipped. In short, lossless LAVQ can achieve rates comparable to the best LZ-based algorithms.

VI. Software Implementation

Software implementation of LAVQ, complete with all improvements, has been completed. To minimize source code complexity, the arithmetic coder has been implemented separately from the encoder. Parameters variable

in the encoder include block dimensions, codebook size, error tolerance, number of bits stripped, first or best occurrence of an acceptable codeword, and difference coding. Furthermore, to allow use of the encoder for image sequences, as in video, codebook preservation is also provided. This allows the codebook from the previous frame to be used in the subsequent one.

A. LAVQ Encoder and Decoder

The encoder program first reads in all the data, then converts the image into a linked list of blocks. The codebook is specified as a doubly linked list to facilitate speed in rearrangement: Only ten pointers at most need be changed to do a complete codebook rearrangement. The codebook is initially assumed empty. At each step, a block is compared with the codebook entries and the best or first occurrence of an acceptable codeword is found. The codebook is updated as needed, and the requisite values are output. The program can also calculate global entropy of the encoded file to give an estimate of the highest compression possible with these parameters. The decoder program reverses this arrangement, with each input being an index or a new codebook value, and the codebook is rearranged in a manner identical to the encoder codebook. The image is rebuilt as a linked list and is then converted to image format and output.

B. Adaptive Arithmetic Coding

The arithmetic coder implemented is one described in [5]; it is used here with only minor modifications; the most notable of these is the ability to input symbols of alphabet size less than 256. The assumption is made at first that the symbols are all equiprobable; after each symbol is encoded, its statistics are modified to reflect this. Two encoders and decoders are used: One pair encodes and decodes only the indices from the output of the LAVQ encoder and ignores the new codeword information; the second pair codes and decodes these new codeword values.

The basic arithmetic code operates in the following manner: The symbols are arranged in some order and are assigned regions of size corresponding to their probabilities. These regions span the space from 0 to 1. The coding region is initially defined to be [0,1). When a symbol is coded, the symbol space [0,1) is scaled to fit the coding region, and the new coding region is defined by the region specified by the symbol coded. Therefore, at each step, the coding region becomes more and more narrow. In practice, the coding region is scaled in size as each unambiguous most significant bit is transmitted. The decoder reverses this by scaling up the coding region as symbols

are decoded. A detailed discussion of arithmetic codes is available in [5].

VII. Experimental Results

A. Image Compression

The LAVQ algorithm was tested on a number of monochrome 8-bit 512×512 pixel images. Global pixel entropies, which are entropies estimated over the whole images, are listed in Table 1. These images were selected to provide a diverse cross section of images to examine the flexibility of LAVQ. A portrait ("lena"), a wildlife/natural scene of a seal on a rocky seashore ("seal"), a high detail overhead view of Los Angeles International Airport ("lax"), the cratered surface of Mercury ("mercury"), the rings of Saturn ("saturn"), and a medical CAT scan image ("cat01") were used. The images are shown in Fig. 2.

Performance parameters are measured in mean squared error (MSE) for distortion and required bits per pixel for rate. MSE is defined as

$$MSE = \frac{1}{512^2} \sum_{i=1}^{512^2} [p_{\text{original},i} - p_{\text{processed},i}]^2$$

for images of size 512×512 pixels.

The LAVQ parameters used were 8×1 blocks with 255 codewords in the codebook. Difference coding was used, and 4 bits with logarithmic quantization levels were used to represent each new codeword value. There is not much difference in performance between finding the best match in the codebook (complete codeword search) and stopping after the first instance of an acceptable (within error allowance) codeword in the book (partial codebook search, typically one-fifth to one-tenth of the entire book); therefore, the latter strategy is used to maximize speed and to favorably skew the index statistics. Codebook and block sizes were selected to obtain good results; block sizes larger or smaller than 8×1 yielded worse results, and increasing codebook size beyond 255 had only marginal rate improvement with substantial increase in computational complexity. Block interpolation and smoothing are used with a threshold of 32. The LBG algorithm using blocks of size 4×4 pixels (found to yield good results in general and better results than LBG with 8×1 , 4×2 , or 16×1 pixel blocks) is presented for comparison.

In the rate-distortion curves generated for all six images (see Figs. 3 through 8), the LAVQ blocks are 8×1 pixels;

255 codewords are in the codebook, the difference coding has 4 bits per codeword value, and there are logarithmic quantization, block interpolation, and smoothing, except for Fig. 5, which lacks block interpolation and smoothing. LBG was done with 4×4 pixel blocks. LAVQ sends the codebook simultaneously with the indices while compressing the image; therefore, a true comparison should include the cost of sending the codebook for LBG as well. The LBG curves were generated using 4×4 pixel blocks and varying codebook sizes (from 16 to 8192). LAVQ curves were generated using the parameters outlined above. The curves have differing scales because each image has different characteristics which alter the algorithms' ability to compress them.

In most cases, LAVQ does better (defined as having a lower distortion for a given rate or vice versa) than LBG in the low-distortion regions (high-rate regions); in this region, LAVQ sends more new codewords, and these new codewords are more accurate renditions of the original image than the codewords used by LBG, which are the centroids of many blocks. Because LAVQ requires that more codewords be sent than LBG for low-distortion cases, this factor by itself would seem to make the rate for LAVQ worse than for LBG. However, LAVQ can take good advantage of lossless entropy coding of indices and new codewords. In contrast, because LBG's codebook search algorithm distributes its codewords to span the space in which the image's vectors exist, in general LBG does not profit as much as LAVQ does from lossless compression. Only about a 5-percent improvement using arithmetic coding was noted, compared to over 50-percent improvement obtained for LAVQ using arithmetic coding of both indices and codewords. As a result, LAVQ often has the potential for better performance than LBG after entropy coding is applied to LAVQ.

At lower rates (and higher distortions), LBG takes time to select a good codebook and therefore can do much better than LAVQ, which does not train a codebook at all. This is the dominant factor which can make LAVQ performance worse than LBG: Since LAVQ does not train a codebook, it relies on recently occurring vectors for a source of new codewords. These do not necessarily reflect a good choice of codewords, so for a given rate, the distortion for LAVQ can be higher than for LBG.

Such generalizations, however, have many exceptions. For example, two images in the set used here, "cat01" and "saturn," have large regions of uniformity compared to the other images. In the high-rate (low-distortion) region, these low-detail regions are better coded by LBG than by LAVQ because the LBG codewords can code regions of

low detail with less distortion; this effect is enough in these two images to defeat the advantage that LAVQ gains by lossless compression of the indices and codewords.

In the low-rate (high-distortion) region, "cat01" and "saturn" are better coded with LAVQ than with LBG. With other images, the time spent by LBG in developing a good codebook made the performance of LBG better than what LAVQ could achieve without training. With these two images, however, LAVQ can code the low-detail regions quite easily with very few new codewords while the details are better preserved in the high-detail regions. This can be enough to offset the disadvantages mentioned earlier that LAVQ encounters in the low-rate (high-distortion) region.

Each LBG rate-distortion curve of 10 data points took approximately 100 hours of computation on a Sun Sparc 2, while one LAVQ curve with 24 data points, complete with lossless compression, was done in about 1 hour: LAVQ is much faster than LBG and can achieve performance comparable to LBG if the cost of the codebook is included. Furthermore, LAVQ preserves detail better than LBG, but does so with increased blockiness in low-detail regions. This detail-preserving feature of LAVQ is discussed in the next section.

B. Detail Preservation

As mentioned in the previous sections, LAVQ operates by matching a vector to the codewords in the codebook using a predetermined fidelity criterion, and new codewords are entered verbatim from the examined vector if no match occurs within this error allowance. Therefore, those vectors which are significantly different from previous vectors are coded without distortion; these occur at edges or areas of high detail. The cost of preserving these details, however, is increased distortion in low-detail ("smooth") areas. Codewords are not optimized to best represent these regions, so they exhibit more blockiness. Thus, LAVQ preserves details and is potentially very attractive to military intelligence and space applications: These applications require close examination of details to identify and differentiate among various objects.

These effects are illustrated in Figs. 9 and 10. In Fig. 9, the upper right edge of the hat brim of "lena" is shown. The LBG codebook size, to be comparable to LAVQ, is fixed at 256. The LAVQ error allowance was adjusted to yield a distortion comparable to LBG. LBG achieves 43.83 MSE at 0.50 bits/pixel (compression ratio of 16:1). LAVQ achieves 43.68 MSE at 0.56 bits/pixel (compression ratio of 14:1). The edge of the hat brim is rendered with much

less blockiness with LAVQ; the drawback is that the low-detail areas exhibit noticeable blockiness.

In Fig. 10, the lower right edge of the upper terminal of "lax" is shown. The LBG codebook size is again fixed at 256. LBG achieves 166.8 MSE at 0.50 bits/pixel (compression ratio of 16:1). LAVQ achieves 166.0 MSE at 0.77 bits/pixel (compression ratio of 10:1). Here, LBG does poorly in preserving the details of the aircraft, terminal, and service vehicles; details of the terminal and several service vehicles have disappeared. LAVQ does much better on these details, but exhibits more blockiness than LBG when representing the tarmac, which has less detail. With LAVQ, it is possible to identify aircraft type, while it is more difficult with LBG. For the two aircraft in Fig. 10, the fuselage and wing shapes, engine locations, and other details are preserved more clearly by LAVQ than by LBG.

C. Lossless Data Compression

Demonstration of lossless compression of structured data was conducted on Magellan space-probe engineering data. The Magellan spacecraft engineering data file "mdata" consists of 1692 records of 100 bytes each; each record consists of 15 fields of various sizes ranging from 4 bytes to 16 bytes. Many good existing LZ variants can achieve compression ratios for "mdata" of 2.2:1 to 5:1 of the original file size, depending on the amount of memory consumed. The UNIX program "compress" achieves a compression of 3:1 at a cost of 544 kbytes of memory used.

The Magellan engineering data have both natural parsing (a frame marker at the beginning of each record) and artificial parsing (the record size is fixed, 100 bytes) properties. Most real-world data sources like texts, images, and engineering data records possess either or both of these two parsing properties. To apply LAVQ to these data, the file was first blocked in sequential order into 5-byte blocks. LAVQ was then applied to each block position; thus, all of the first 5-byte blocks (bytes 1 to 5 in the 100-byte record) were coded using one encoder, all of the second 5-byte blocks (bytes 6 to 10 in the 100-byte record) were coded using a second encoder, and so on.

The resulting LAVQ output, both indices and new codeword values, was then coded with the adaptive arithmetic coder used earlier. Results using a Q-coder [6,7] and the theoretical limit reached by nonadaptive means (computed from the global entropy) are also included for comparison. The Q-coder is essentially an adaptive arithmetic coder; it differs from the arithmetic coder used above: The arithmetic coder maintains a running total of frequency of use of each symbol and uses this to determine probabilities

for arithmetic coding. The Q-coder uses a fixed probability table accessed by an adaptive finite state machine. This finite state machine adapts with the previous symbols; its state is determined by the bits in each byte.

Two differing approaches are taken. First, the codebook size is fixed at 256, and only one coder is available. Thus, a large buffer is needed to store the data to allow sequential coding of each of the 20 blocks in each record. Results tabulated in Table 2 show that the best performance is obtained from using large buffers: Larger files have more opportunities for repetitions and patterns of codeword usage. In Table 2, compression ratios of Magellan engineering data are of 169,200-byte size. Tabulated are compression ratios achieved by the adaptive arithmetic coder used for LAVQ, the Q-coder, and the theoretical maximum using a global entropy coder. This last value is derived from the global per-symbol entropy. Note that larger buffer sizes provide greater compression. Arithmetic coder compression approaches the theoretical value (global entropy) closely. Even better results are sometimes obtained. The Q-coder and arithmetic coder sometimes appear to have results better than entropy; this occurs because these adaptive coders code on the basis of local entropy (entropy over a localized window of characters) as opposed to the global entropy listed here. The Q-coder does better than arithmetic coding, which does about as well as theoretically possible with a global entropy coder.

In the second approach, the codebook size is allowed to vary, and all 20 coders are available for parallel coding. No additional buffer is needed in this case, as the 20 blocks are coded simultaneously. Again, the indices and new codeword values are coded with the adaptive arithmetic coder and compared with the Q-coder and the global entropies; results are listed in Table 3. In this table, as in Table 2, compression ratios of Magellan engineering data are of 169,200-byte size. No data buffer is allowed here, but 20 coders are operating in parallel. In this case, the arithmetic coder does slightly better than the Q-coder. Again, the Q-coder and arithmetic coder sometimes appear to have results better than entropy for the same reason given for Table 2. Best performance is obtained with larger codebooks; large codebooks can record codewords farther into the past and therefore have more opportunities for codeword matching. Here, the arithmetic coder does better than entropy and the Q-coder. In both cases, LAVQ performance is comparable to LZ-based algorithms. The advantage of LAVQ is that far less memory is required than for LZ algorithms; this is of importance in systems with size, weight, or complexity constraints, as is the case with deep-space probes.

VIII. Conclusion

The LAVQ algorithm provides a fast, one-pass data compression algorithm. Improvements to the basic algorithm maintain this one-pass high-speed property while increasing performance measurably. Experimental results

in image compression yield performance not significantly inferior to LBG, but at a fraction of the complexity. Distortion in images occurs as blockiness in low-detail areas, while high-activity areas maintain sharp details. LAVQ also performs well in lossless compression, again with low complexity as compared with other algorithms.

References

- [1] Y. Linde, A. Buzo, and R. M. Gray, "An Algorithm for Vector Quantization," *IEEE Transactions on Communications*, vol. C-28, pp. 84-89, 1980.
- [2] J. L. Bentley, D. D. Sleator, R. E. Tarjan, and V. K. Wei, "A Locally Adaptive Data Compression Scheme," *Communications of the ACM*, vol. 29, pp. 320-330, 1986.
- [3] K. M. Cheung and V. K. Wei, "A Locally Adaptive Source Coding Scheme," *Proceedings of the Bilkent Conference on New Trends in Communications, Control, and Signal Processing*, Ankara, Turkey, pp. 1473-1482, July 2-5, 1990.
- [4] K. M. Cheung and V. K. Wei, "A Locally Adaptive Source Coding Scheme," submitted to *IEEE Transactions on Communications*, 1990.
- [5] I. H. Witten, R. M. Neal, and J. G. Cleary, "Arithmetic Coding for Data Compression," *Communications of the ACM*, vol. 30, pp. 520-540, 1987.
- [6] Joint Photographic Experts Group, *JPEG Draft Technical Specification*, Rev. 5, Washington, DC: International Organization for Standardization/International Telegraph and Telephone Consultative Committee (ISO/CCITT), Section 12, January 15, 1990.
- [7] R. B. Arps, T. K. Truong, D. J. Lu, R. C. Pasco, and T. D. Friedman, "A Multi-Purpose VLSI Chip for Adaptive Data Compression of Bilevel Images," *IBM Journal of Research and Development*, vol. 32, pp. 775-795, 1988.

Table 1. Pixel entropies. Global pixel entropies of images.

Images	Entropy, bits/pixel
"cat01"	5.503
"lax"	6.827
"lena"	7.445
"mercury"	6.416
"saturn"	6.885
"seal"	7.356

Table 2. Magellan data compression: sequential compression.

Buffer size, bytes	Codebook size	Compression ratio achieved by		
		Arithmetic coder	Q-coder	Global entropy
169200	256	3.957	5.053	3.887
84600	256	3.752	4.764	3.748
56400	256	3.635	4.556	3.640
42300	256	3.429	4.210	3.440
28200	256	3.285	3.929	3.296
14100	256	2.913	3.320	2.921

Table 3. Magellan data compression: simultaneous compression.

Codebook size	Compression ratio achieved by		
	Arithmetic coder	Q-coder	Global entropy
256	3.957	3.912	3.887
128	3.873	3.815	3.798
64	3.411	3.290	3.314
32	2.973	2.814	2.949
16	2.747	2.545	2.6076

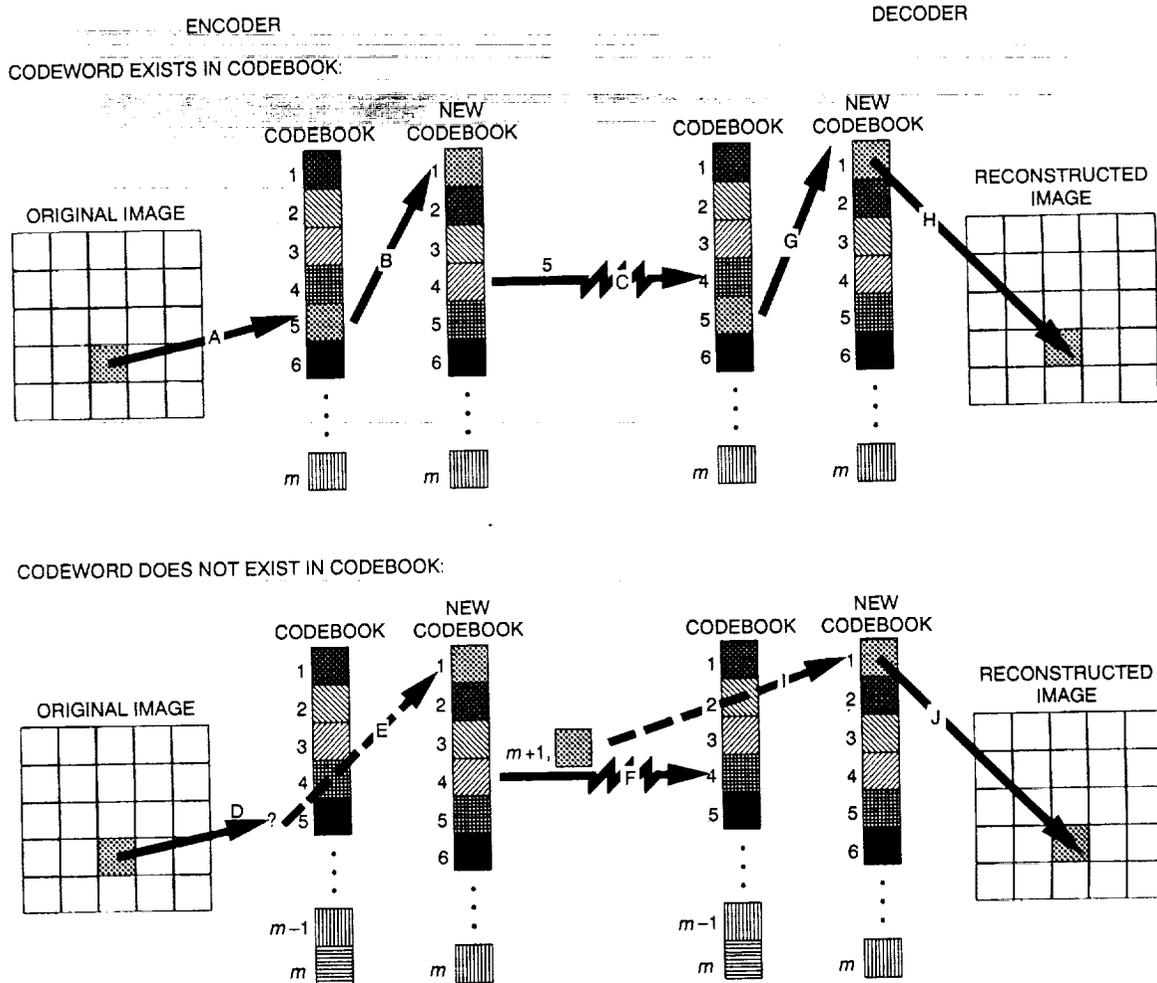
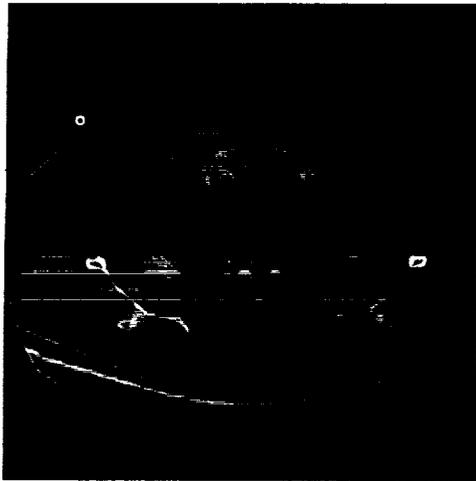
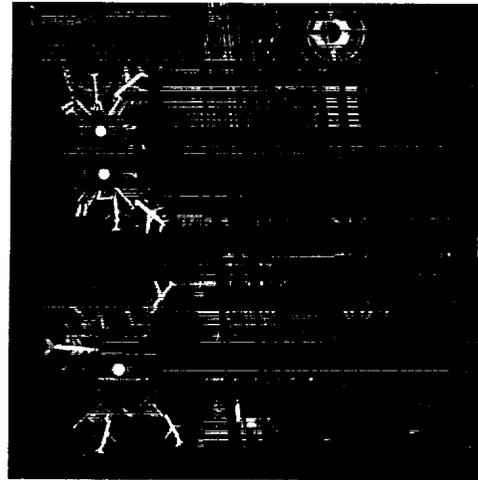


Fig. 1. Encoder and decoder: An image block is compared to the codebook (A); if a codeword close enough exists, that codeword is moved to the top of the codebook (B) and the index is transmitted (C). If it does not exist (D), then the block is inserted at the top of the codebook (E) and the index $m + 1$ and the block are transmitted (F). On the receiver side, if an index is received, the corresponding codeword (G) is moved to the top of the codebook and the block is inserted into the reconstructed image (H). If the special index $m + 1$ is received (I), a raw block is anticipated immediately following; this block is placed in the codebook and also in the reconstructed image (J).



"cat01"



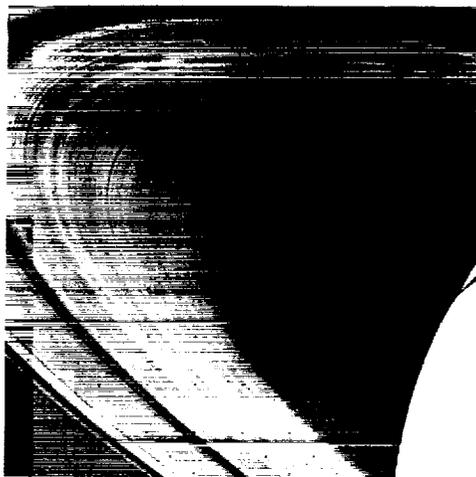
"lax"



"lena"



"mercury"



"saturn"



"seal"

Fig. 2. Original Images. All images are 512 × 512 pixel, 8-bit monochrome.

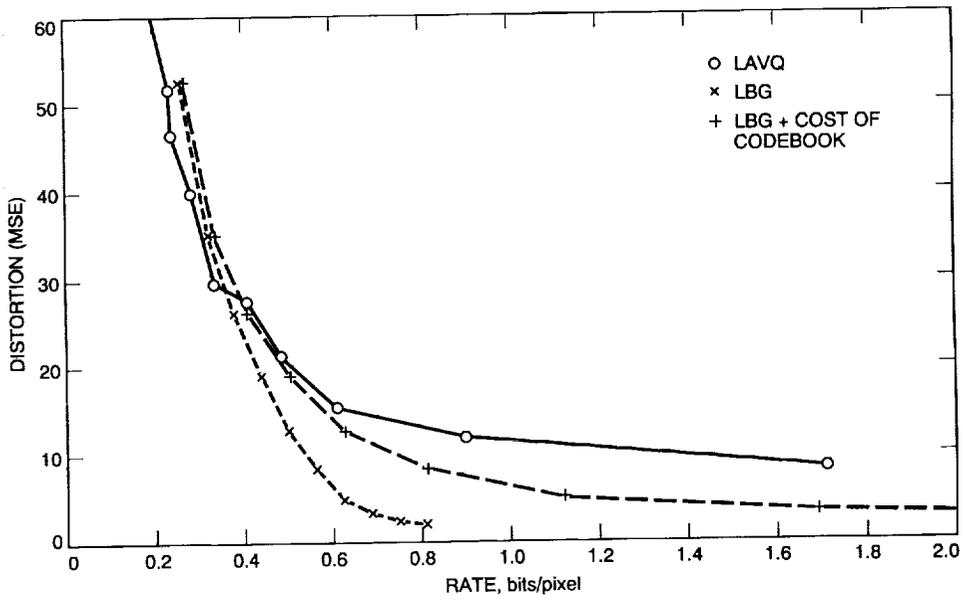


Fig. 3. Rate-distortion curve for "cat01."

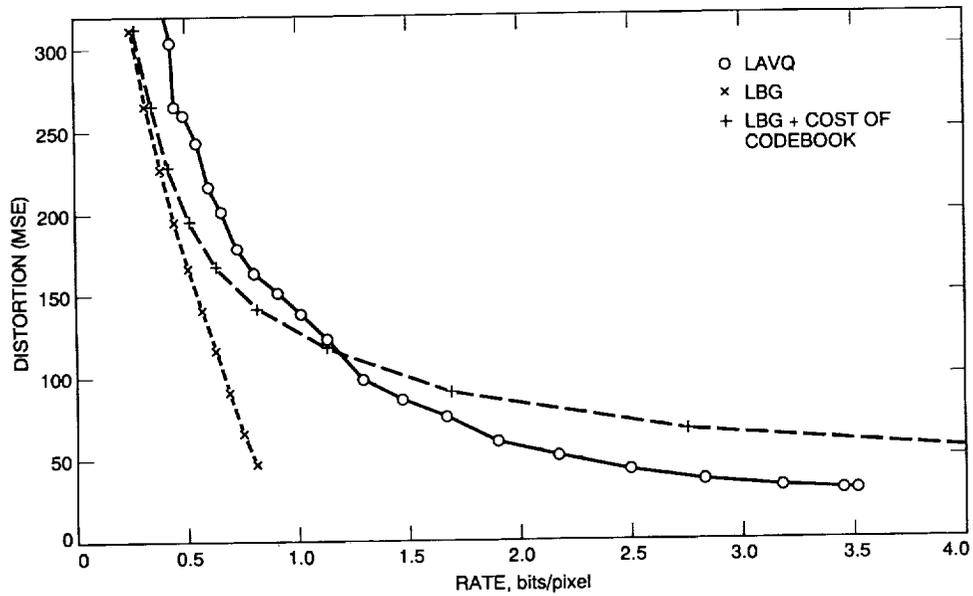


Fig. 4. Rate-distortion curve for "lax."

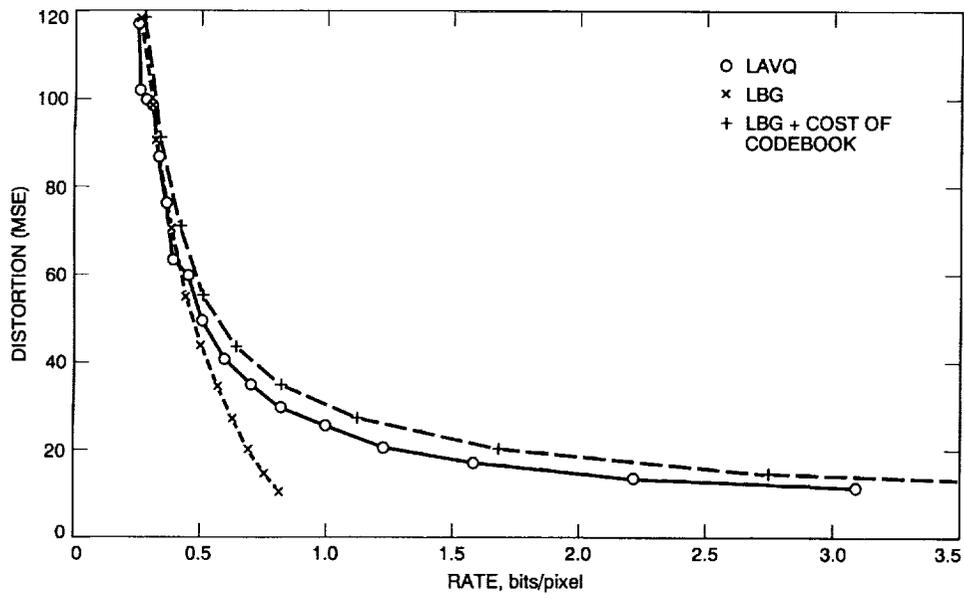


Fig. 5. Rate-distortion curve for "lena."

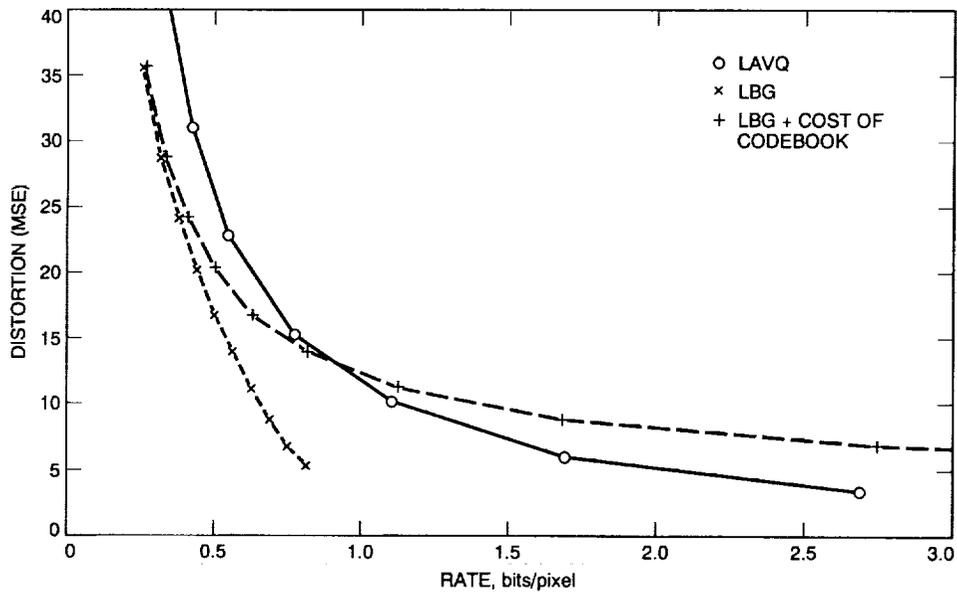


Fig. 6. Rate-distortion curve for "mercury."

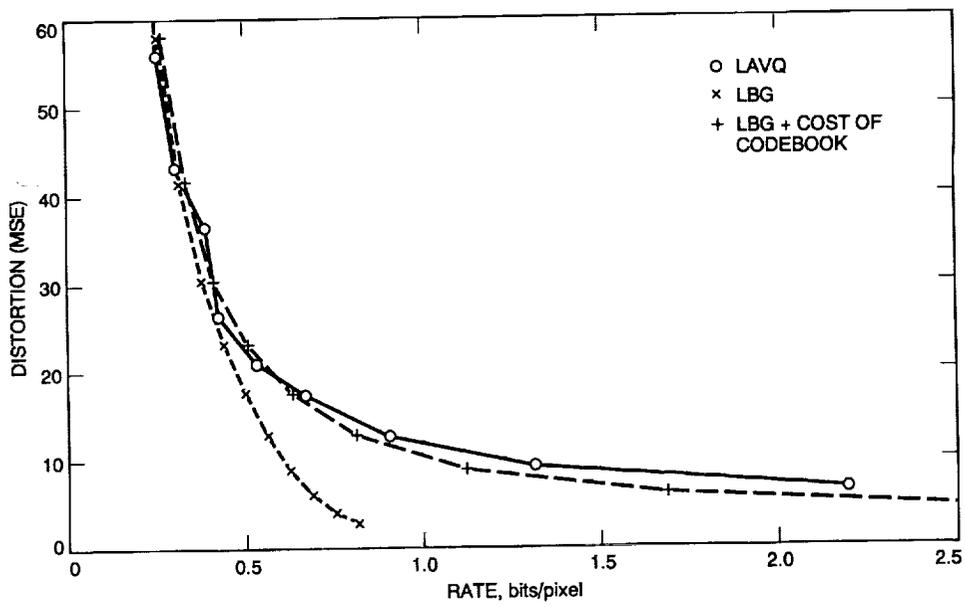


Fig. 7. Rate-distortion curve for "saturn."

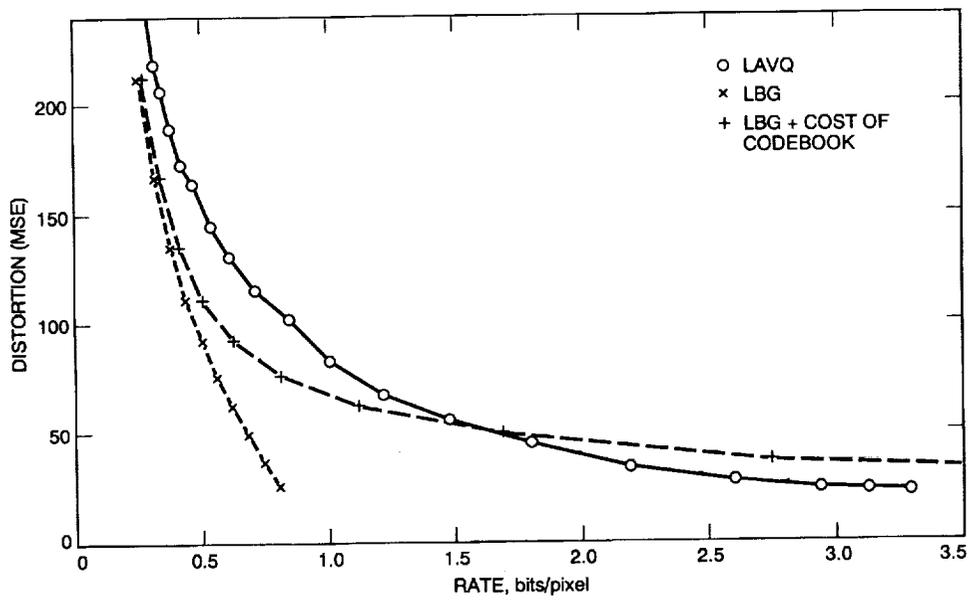
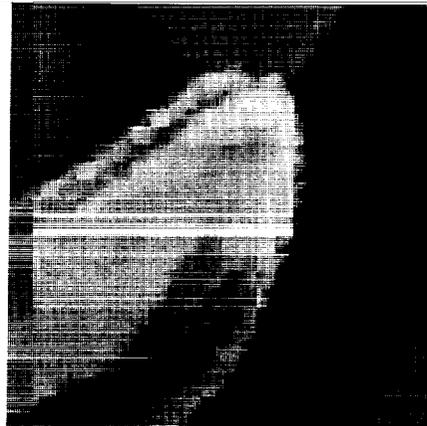


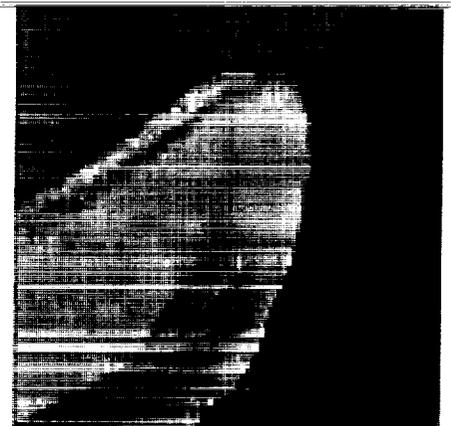
Fig. 8. Rate-distortion curve for "seal."



ORIGINAL

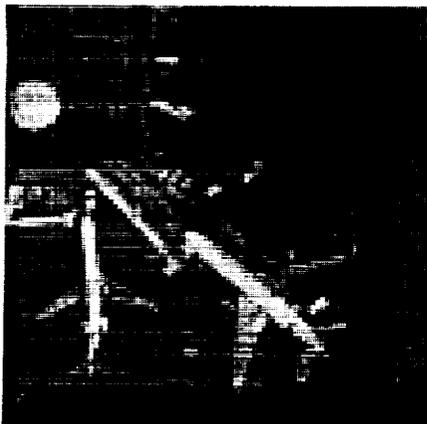


LBG

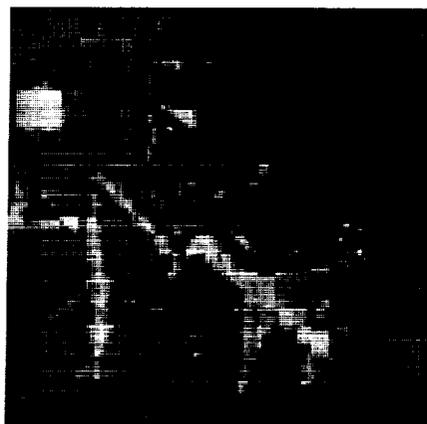


LAVQ

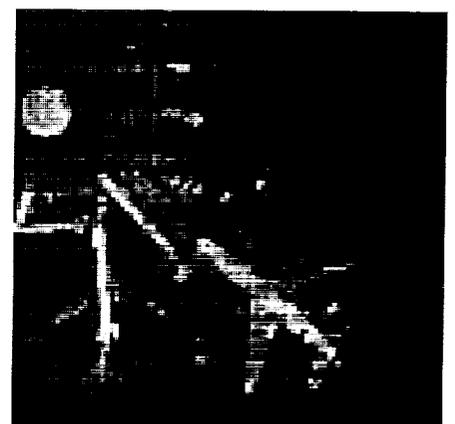
Fig. 9. Detail of "lena." Note that LBG has more blockiness at the edge, but represents low-detail ("smooth") areas without as much blockiness as LAVQ.



ORIGINAL



LBG



LAVQ

Fig. 10. Detail of "lax."

513-61
N93-19426
R8446
p-9

Data Compression by Wavelet Transforms

M. Shahshahani

Communications Systems Research Section

A wavelet transform algorithm is applied to image compression. It is observed that the algorithm does not suffer from the blockiness characteristic of the DCT-based algorithms at compression ratios exceeding 25:1, but the edges do not appear as sharp as they do with the latter method. Some suggestions for the improved performance of the wavelet transform method are presented.

I. Introduction

The application of wavelet transforms and multiresolution analysis to data compression has attracted much attention recently. This circle of ideas is closely related to the subband compression and the pyramid encoding techniques. The general idea is to transform and reorganize the data in a hierarchical manner so that the upper levels of this hierarchy (or pyramid) represent the general features of the data or the image and the lower levels supply the details. Generally the higher levels of the pyramid are smaller data sets than the lower levels; however, the coefficients in the latter portion are more correlated than those in the former and are better compressed by the standard lossless compression techniques.

The applications of wavelet representations to practical engineering problems are not limited to source coding. For example, one encounters situations that necessitate selecting a subset of a large data set on the basis of certain characteristics. One may achieve this by browsing through the higher levels of the hierarchy, which comprise a much smaller data set, examining the general features of

the data, and making judicious choices. The coefficients in the lower levels of the pyramid may be used for edge detection.

As in other methods of data compression, applications of wavelet transforms to source coding assume a priori knowledge of the tolerable level of information loss and/or the desirable compression ratio. Data compression is achieved by quantizing the transformed data and allocating bits to the different levels of the pyramid of the transformed data in a manner compatible with the constraints and the requirements of the particular application. Naturally, in source coding applications more bits are allocated to an individual coefficient in the higher levels of the pyramid than to one in a lower level. In analogy with the discrete cosine transform (DCT), one may regard the lower levels of the pyramid as the high frequencies and the upper ones as the low frequencies.

The presentation of the general theory of wavelet transforms in Section II is intended for application to data compression. The literature on the subject is often inadequate regarding the implementation of the basic ideas, and the

theoretical aspects of the subject seem to be only remotely related to practical engineering problems. It is hoped that the concise and concrete presentation of the wavelet transforms in Section II will make the literature more accessible to interested researchers. In Section III, the practical aspects of image compression by wavelet transforms and the results of the applications are reported. Further research topics for the improvement of the performance of wavelet-transform-based compression algorithms are also suggested. Some of the advantages and disadvantages of the wavelet transforms versus the standard DCT techniques are discussed. However, no definitive judgment can be made at this time regarding their relative merits. While the latter approach has been studied extensively in the past decade, the application of wavelet transforms to image compression has not reached the level of maturity that would warrant definitive assessment of its merits and potential.

II. Wavelet Transforms

The idea of wavelet transforms and their applicability to signal analysis, and especially data compression, is most easily demonstrated by focusing on the one-dimensional case first. A straightforward generalization of the theory to two dimensions for application to image compression is indicated at the end of this section. In this case, a data set is represented by an element of $\mathcal{L} = L^2(\mathbf{R})$. Consider the following sequence of partitions of \mathbf{R} :

$$\text{Partition } \mathcal{P}_m : \quad \mathbf{R} = \bigcup_{n=-\infty}^{\infty} I_{n,m}$$

where $I_{n,m} = [2^m n, 2^m(n+1))$, and let \mathcal{L}_m be the subspace of \mathcal{L} consisting of functions that are constant on the intervals $I_{n,m}$. The operator p_m of orthogonal projection on the subspace \mathcal{L}_m is

$$p_m(f)(x) = \frac{1}{2^m} \int_{I_m(x)} f(x) dx$$

where $I_m(x)$ is the unique interval $I_{n,m}$ (m fixed) containing x . The subspaces \mathcal{L}_m have the following properties:

$$\mathcal{L}_{m+1} \subseteq \mathcal{L}_m, \quad \bigcap \mathcal{L}_m = 0, \quad \overline{\bigcup \mathcal{L}_m} = \mathcal{L} \quad (1)$$

Let \mathcal{E}_m denote the orthogonal complement of \mathcal{L}_{m+1} in \mathcal{L}_m , then \mathcal{L} admits of the orthogonal direct sum decomposition $\mathcal{L} = \bigoplus \mathcal{E}_m$. Denote orthogonal projection on \mathcal{E}_m by π_m . Let $A_{a,b}$, where $a \neq 0$ and b are real numbers, denote the affine transformation $A_{a,b}(x) = ax + b$, and de-

fine the action of $A_{a,b}$ on a function φ by $A_{a,b}(\varphi)(x) = a^{-1/2} \varphi[(x-b)/a]$. It is convenient to introduce the notation $\varphi_{m,n}(x) = A_{2^m, 2^m n}(\varphi)(x)$, for m and n integers, and note that a function $f \in \mathcal{L}$ admits of the expansion

$$p_m(f) = \sum_{n=-\infty}^{\infty} a_n^m \chi_{m,n} \quad (2)$$

where χ is the indicator function of the interval $[0,1)$, and $a_n^m = \int \chi_{m,n}(x) f(x) dx$. The functions $\chi_{m,n}$ are obtained from the single function χ through the action of a set of affine transformations of the line. For each fixed m ,

$$\mathcal{L}_m = \text{span} \{ \chi_{m,n} | n \in \mathbf{Z} \} \quad (3)$$

$$f \in \mathcal{L}_m \iff f(2 \cdot) \in \mathcal{L}_{m-1}$$

From the expansion (2) one easily obtains the expansion of f following the decomposition $\mathcal{L} = \bigoplus \mathcal{E}_m$. First observe that

$$\chi_{m+1,n} = \frac{1}{\sqrt{2}} (\chi_{m,2n} + \chi_{m,2n+1})$$

Therefore, $a_n^{m+1} = \frac{1}{\sqrt{2}} (a_{2n}^m + a_{2n+1}^m)$, and after a simple calculation one obtains

$$p_m(f) - p_{m+1}(f) = \frac{1}{2} \sum (a_{2n}^m - a_{2n+1}^m) (\chi_{m,2n} - \chi_{m,2n+1})$$

Now set $\varphi_{m,n} = \frac{1}{\sqrt{2}} (\chi_{m,2n} - \chi_{m,2n+1})$ to obtain the expansion

$$f = \sum_{m,n=-\infty}^{\infty} b_n^m \varphi_{m,n} \quad (4)$$

where $b_n^m = \frac{1}{\sqrt{2}} (a_{2n}^m - a_{2n+1}^m)$. It is a remarkable fact, and easy to prove, that the functions $\varphi_{m,n}$ are also obtained from the single function $\varphi(x) = \chi(2x) - \chi(2x-1)$ by the action of the set $\mathcal{A} = \{ A_{2^m, 2^m n} | m, n \in \mathbf{Z} \}$ of affine transformations, and an analogue of condition (3) is valid for the subspaces \mathcal{E}_m , namely,

$$\mathcal{E}_m = \text{span} \{ \varphi_{m,n} | n \in \mathbf{Z} \} \quad (5)$$

$$f \in \mathcal{E}_m \iff f(2 \cdot) \in \mathcal{E}_{m-1}$$

The functions $\{\varphi_{m,n}\}$ form a complete orthonormal set for \mathcal{L} . The expansion (4) is an example of orthonormal wavelet expansion, and the coefficients b_n^m are called the wavelet coefficients.

To understand the intuitive meaning of the expansion (4), assume that $f \in \mathcal{L}_m$ for a sufficiently large negative number m . The projection of f on \mathcal{E}_m is $\pi_m(f) = f - p_{m+1}(f)$. Now $p_{m+1}(f)$ is a slightly *smoothed* version of f so that $\pi_m(f)$ represents the details that are missing from the smoothed version $p_{m+1}(f)$. Thus, $\pi_m(f)$ or more precisely, the coefficients b_n^m in the expansion (4) belong to the lowest level of the hierarchy. The process can be repeated with $p_{m+1}(f)$ replacing f , thus leading to a hierarchy of the coefficients of the wavelet expansion of f .

An important feature of the expansion (4) is that the coefficients b_n^m and a_n^m can be computed recursively in a simple manner. For example, to compute a_n^1 from a_n^0 , substitute expansion (2) for $m = 0$ in the formula for a_n^1 to obtain

$$a_n^1 = \sum_{k=-\infty}^{\infty} \int \frac{1}{\sqrt{2}} \chi(x-k-2n) \chi(x/2) dx \quad (6)$$

Therefore, if one defines $\alpha(k)$ as the integral in expansion (6) for $n = 0$, one obtains the formula

$$a_n^1 = \sum_{k=-\infty}^{\infty} \alpha(k-2n) a_k^0 \quad (7)$$

Similarly,

$$b_n^1 = \sum_{k=-\infty}^{\infty} \beta(k-2n) a_k^0 \quad (8)$$

where $\beta(k) = \frac{1}{\sqrt{2}} \int \chi(x-k) \varphi(x/2) dx$. By a straightforward inductive extension of this calculation, one can express b_n^{m+1} and a_n^{m+1} in terms of a_n^m . The resulting formulae are identical with formulae (7) and (8) with m and $m+1$ replacing 0 and 1, respectively. Therefore, the wavelet coefficients b_n^m and a_n^m can be computed by the filters defined by α and β .

The orthonormal basis $\{\varphi_{m,n}\}$ and the expansion (4) are just one example of an orthonormal wavelet expansion. To obtain other expansions, one has to abstract some of the features of this illustrative example. The essential ingredients of the theory are an orthonormal doubly infinite

basis $\{\varphi_{m,n}\}$ for \mathcal{L} such that the functions $\varphi_{m,n}$ are obtained from a single function φ via the action of the set \mathcal{A} , and for which condition (5) is valid. For applications, knowledge of the corresponding filters β and α is essential. Since in practical engineering problems the data are normally in digital form, it is important to adapt the theoretical framework of wavelets to the discrete or digital case before discussing other wavelet expansions.

In the digital domain, $l^2(\mathbf{Z})$ replaces $L^2(\mathbf{R})$ as the space of one-dimensional data. One can naturally identify $l^2(\mathbf{Z})$ with \mathcal{L}_0 , and therefore the theory developed above extends to this case immediately. The only difference is

$$\mathcal{L}_{-j} = \mathcal{L}_0, \quad p_{-j} = id., \quad \text{and} \quad \pi_{-j} = 0 \quad \text{for } j \geq 0 \quad (9)$$

It follows that formulae (2) through (8) remain valid, provided that the range of the values of m is limited to 0 to ∞ . In practice, the domain of n is $(\mathbf{Z} \bmod 2^N)$ for some integer N . Therefore, $\mathcal{L}_N = \mathbf{R}$, and the linear spaces \mathcal{L}_m are finite dimensional. The bases $\{\chi_n^m\}$ and $\{\chi_n^{m+1}, \varphi_n^{m+1}\}$ for $\mathcal{L}_m = \mathcal{L}_{m+1} \oplus \mathcal{E}_m$ differ by an orthogonal transformation. It follows that the coefficients $\{a_n^m\}$ and $\{a_n^{m+1}, b_n^m\}$ are also related by an orthogonal transformation. This orthogonal transformation, which is the matrix representation of the filters α and β , is given by the $2^N \times 2^N$ matrix with 2×2 diagonal blocks

$$\begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix}$$

This means that given a data set represented by a column vector $(f_0, \dots, f_{2^N-1})^{tr}$, the application of the above matrix transforms it into a vector $(g_0, \dots, g_{2^N-1})^{tr}$ with the even-numbered components $(g_0, g_2, \dots, g_{2^N-2})^{tr}$ representing $p_1(f)$ and the odd-numbered ones $(g_1, g_3, \dots, g_{2^N-1})^{tr}$ representing $\pi_1(f)$. Here the superscript tr means the transpose of the matrix or vector.

The problem of determining other orthonormal wavelet expansions, and especially the corresponding filters, is discussed in detail in [1]. Of particular interest in practical problems is the case where the functions α and β have small support, i.e., $\alpha(j) = 0 = \beta(j)$ for most j 's. It is the knowledge of the functions (or filters) α and β , and not the basis functions themselves, that is essential for applications. In [1], the filters α and β of small support are explicitly determined. The simplest of these filters is the one given above. The next simplest one is the matrix \mathcal{F} given by

$$\mathcal{F} = \begin{pmatrix} \alpha_0 & \alpha_1 & \alpha_2 & \alpha_3 & 0 & 0 & 0 & \dots & \dots & 0 \\ \alpha_3 & -\alpha_2 & \alpha_1 & -\alpha_0 & 0 & 0 & 0 & \dots & \dots & 0 \\ 0 & 0 & \alpha_0 & \alpha_1 & \alpha_2 & \alpha_3 & 0 & \dots & \dots & 0 \\ 0 & 0 & \alpha_3 & -\alpha_2 & \alpha_1 & -\alpha_0 & 0 & \dots & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \ddots & \ddots & & & \vdots \\ \alpha_2 & \alpha_3 & 0 & 0 & \dots & \dots & \dots & 0 & \alpha_0 & \alpha_1 \\ \alpha_1 & -\alpha_0 & 0 & 0 & \dots & \dots & \dots & 0 & \alpha_3 & -\alpha_2 \end{pmatrix} \quad (10)$$

where

$$\alpha_0 = \frac{1 + \sqrt{3}}{4\sqrt{2}}, \alpha_1 = \frac{3 + \sqrt{3}}{4\sqrt{2}}, \alpha_2 = \frac{3 - \sqrt{3}}{4\sqrt{2}}, \alpha_3 = \frac{1 - \sqrt{3}}{4\sqrt{2}}$$

One should note that an important feature of the orthonormal wavelet expansion is that the inversion procedure can be implemented by the transpose of the orthogonal matrix representing the filters α and β .

The above theory was limited to one-dimensional data. It can be easily adapted to the two-dimensional case by considering products of the basis functions considered in the one-dimensional case. This is equivalent to carrying out the one-dimensional wavelet transforms in the horizontal and vertical directions. The practical aspects of the two-dimensional wavelet transform are discussed in detail in the next section. Of course, there are orthonormal wavelet expansions that may not be separable, i.e., the basis functions are not products of the basis functions for the one-dimensional case, but they will not be considered in this article.

III. Application to Data Compression

To apply the theory to data compression, one fixes an orthonormal wavelet expansion, or equivalently, the filters α and β . In the work reported here, only the filter defined by the matrix \mathcal{F} was used.¹ An image is represented by a matrix $f = (f_{ij})$, where f_{ij} is the intensity of the pixel (i, j) . For each fixed row i , one considers the transform $g_i = F f_i^t$, where f_i is the i th row of the matrix f . The components of g_i with even indices represent $p_1(f_i)$ and those with odd indices represent $\pi_1(f_i)$. It is convenient to reorganize the vector g_i in the form

¹ It is often unclear from the literature what filter is actually used. The filter used in [3] differs from that defined by \mathcal{F} .

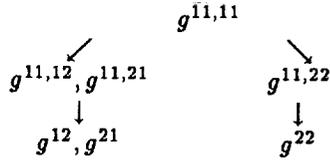
$(g_0, g_2, \dots, g_{2^{N-2}}, g_1, g_3, \dots, g_{2^{N-1}})$. Now the process is repeated for the columns of the matrix of the transformed rows. After reorganizing, the transformed matrix of pixel intensities takes the form

$$g = \begin{pmatrix} g^{11} & g^{12} \\ g^{21} & g^{22} \end{pmatrix}$$

where each g^{ij} is a $2^{N-1} \times 2^{N-1}$ matrix. Since an image is two-dimensional, the hierarchy of the wavelet coefficients requires some elaboration. The matrix g^{11} represents the smoothed version of the image, while the remaining coefficients are the missing details. The coefficients g^{12} and g^{21} belong to the level of the pyramid immediately below g^{11} , and g^{22} lies at the lowest level of the pyramid. Thus, every application of the wavelet transform generates three levels of hierarchy for a two-dimensional image. The process is then repeated by applying the filters α and β to the $2^{N-1} \times 2^{N-1}$ matrix g^{11} along rows and columns. The resulting coefficients are then reorganized in the form

$$g^{11,11} \rightarrow \{g^{11,12}, g^{11,21}\} \rightarrow g^{11,22} \rightarrow \{g^{12}, g^{21}\} \rightarrow g^{22}$$

with the highest level at the extreme left and the lowest at the extreme right. The process can be repeated. It may be more convenient to organize the coefficients differently in the following form:



One then refers to $g^{11,11}$ as SS (smooth-smooth) level 2, to $g^{11,12}, g^{11,21}$ and g^{12}, g^{21} as the SD (smooth-detail) levels 2 and 1, respectively, and to $g^{11,22}$ and g^{22} as the DD (detail-detail) levels 2 and 1, respectively.

In the application of wavelet transforms to image compression, the coefficients at different levels of the pyramid are not equally significant and, therefore, should be encoded differently. The wavelet coefficients of different levels were examined for several images, and certain patterns were observed. In general, the coefficients at a lower level of the pyramid are better approximated by a Laplacian density function than those at the higher levels. Using the nearest integer truncation, one also notices that the entropies of the coefficients at the lower levels are smaller

than those in the upper ones. Figures 1 through 4 show the distributions of the wavelet coefficients at different levels of the hierarchy for a typical image. An approximating Laplacian density function is given in Figs. 1 through 3. Clearly the coefficients at the highest level (Fig. 4) have a very irregular distribution. Table 1 shows the entropies of the wavelet coefficients for the same image.

The image compression process is done by first computing the coefficients $g^{22}, g^{12}, g^{21}, g^{11,22}, g^{11,12}$, etc. These coefficients are quantized according to a bit allocation scheme similar to the one used for the standard DCT-based algorithms. As noted above, more bits are allocated to the higher levels of the hierarchy than to the lower ones. In the pictures of the peppers compressed by the wavelet transform method (Fig. 1), the coefficients in the lowest level have been set to 0. In practice it was observed that because of the quantization errors inherent in any floating-point computation, it is not desirable to go beyond three or four levels of wavelet transforms. To reconstruct the image, the inverse filter was applied to the coefficients. As noted above, the inverse filter is given by the transpose of the orthogonal matrix defining the filter.

There are several issues involved in the application of wavelet transforms to image compression. The choice of the appropriate wavelet transform may be dictated by the complexity of the image. It has been suggested that different transforms may be more appropriate for different images or even different parts of an image. Some ideas in this direction appear in [2] with apparently very promising results. The problem of bit allocation and quantization of the wavelet coefficients is similar to the analogous problem for DCT-based image compression. It may be possible to take advantage of the regularity of the coefficients at the lower levels of the pyramid and use the Laplacian distribution to allocate bits accordingly. However, the experimental work carried out by the author suggests that the simpler method of truncation to the nearest integer followed by decimation by an appropriate number of bits provides better results. Naturally, fewer bits are allocated to the lower levels of the pyramid than to the upper levels. A different method for quantization is proposed in [3]. These authors suggest that using the L^1 rather than the L^2 norm is more compatible with the human visual perception, and their proposed technique of quantization

method is based on minimizing the errors in the former norm.

While a definitive comparison between the DCT-based algorithms and wavelet transform techniques is premature, the tests done by the author suggest some important differences. At higher compression ratios, for example at greater than 25:1, the blockiness in the DCT-based techniques becomes very visible. With the wavelet transform used in the tests, the edges were not as clearly defined as those using the DCT-based techniques, but no blockiness was visible. The rms error of the Joint Photographic Experts Group (JPEG) DCT-based algorithm was smaller than that of the wavelet transform method, but visual preference is not necessarily reflected by the mean square error. Figure 5 shows an original image (peppers) on the upper left corner. The images on the upper right and lower left were compressed using the wavelet transform. The compression ratios were 10:1 and 30:1, respectively. The image in the lower right was obtained by the application of the standard DCT-based JPEG algorithm. Its compression ratio is 30:1. The rms error for the lower left image is about 11.0, and for the one at lower right it is approximately 7.2, even though the blockiness makes it much worse than the one at lower left. The rms error for the image on the upper right is about 9.4. It should be pointed out that other methods, such as fractal algorithms, may produce images that are visually preferable to the DCT-based methods for high compression ratios.

The wavelet transform used in this work is the product of a one-dimensional algorithm with itself; that is, essentially separable into one-dimensional algorithms. It is possible to modify this method to make the horizontal and vertical directions more coupled so that the algorithm becomes truly two-dimensional. The visual effects of such modification are unclear at this time. However, it is reasonable to expect improvements in the clarity of the edges if such techniques are properly employed.

IV. Conclusion

The wavelet transform method provides a new approach to image compression. Although this approach has not performed as well as the DCT-based algorithms in terms of the rms error, it appears to have certain visual advantages especially regarding blockiness.

Acknowledgments

The author wishes to thank Dr. A. Zandi of the University of California at Santa Cruz and IBM (Almaden) for helpful discussions, and Dr. F. Pollara of JPL for providing the DCT-based compressed image of peppers.

References

- [1] I. Daubechies, "Orthonormal Bases of Compactly Supported Wavelets," *Communications on Pure and Applied Mathematics*, vol. 41, pp. 909-996, 1988.
- [2] R. R. Coifman and M. V. Wickerhausen, "Entropy-Based Algorithms for Best Basis Selection," *IEEE Transactions on Information Theory*, vol. 38, no. 2, pp. 713-718, March 1992.
- [3] R. A. DeVore, B. Jawerth, and B. J. Lucier, "Image Compression Through Wavelet Transform Coding," *IEEE Transactions on Information Theory*, vol. 38, no. 2, pp. 719-746, March 1992.

Table 1. Entropies of the wavelet coefficients.

Level	Entropy
DD-1	2.5
SD-1	3.2
DD-2	2.8
SD-2	3.8
DD-3	3.4
SD-3	4.5
DD-4	4.0
SD-4	5.0
SS-4	6.4

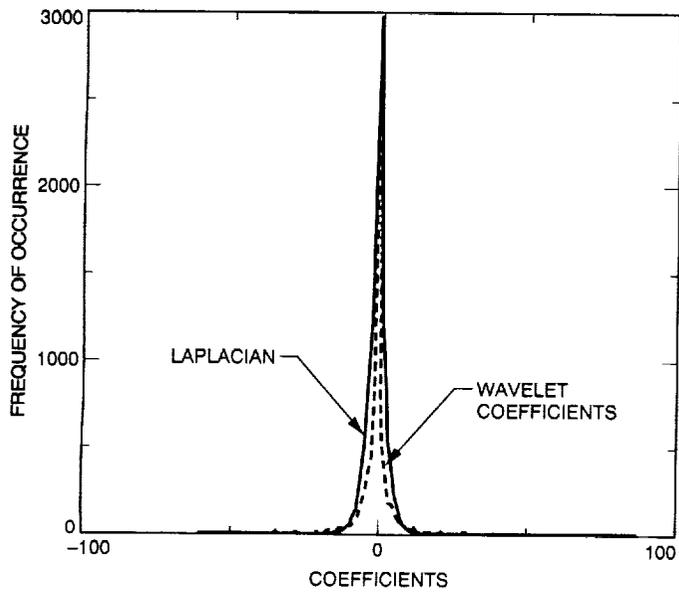


Fig. 1. Level DS-2 wavelet coefficients and Laplacian density.

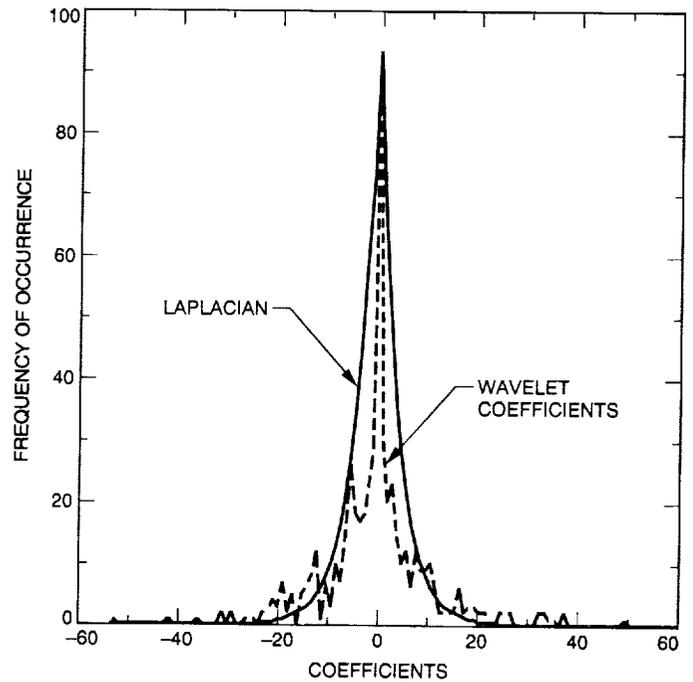


Fig. 3. Level DS-4 wavelet coefficients and Laplacian density.

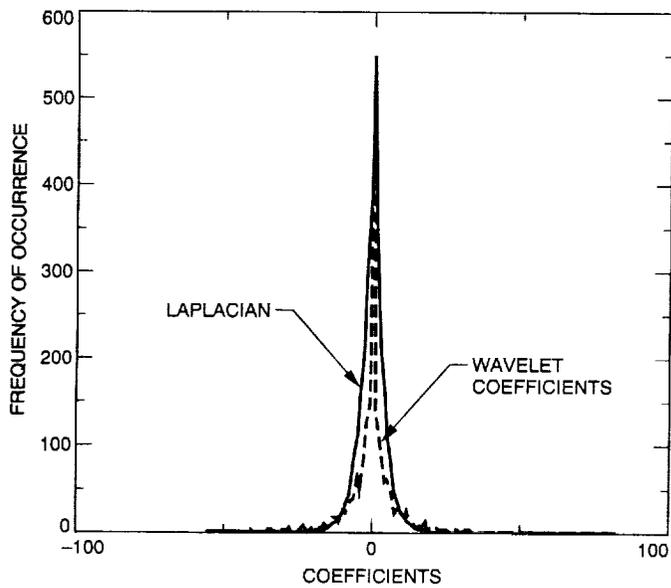


Fig. 2. Level DS-3 wavelet coefficients and Laplacian density.

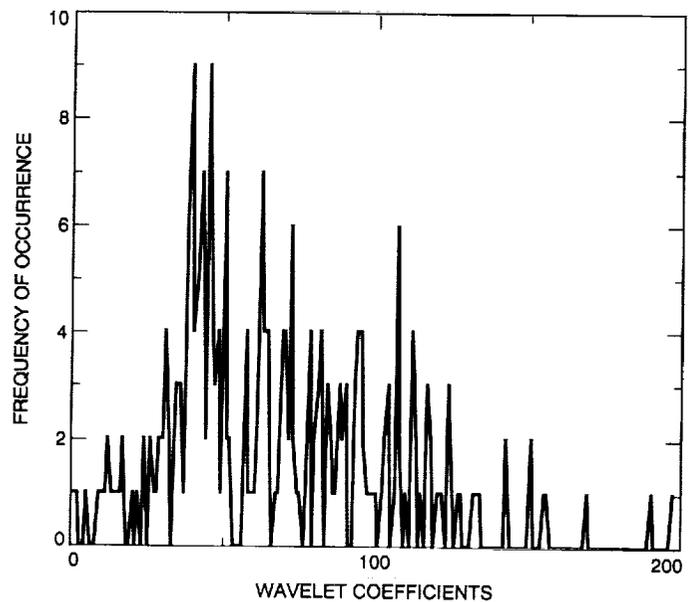


Fig. 4. Level SS-4 wavelet coefficients.



Fig. 5. Image compression: (a) the original uncompressed image; (b) compression ratio of 10:1 by wavelet transform; (c) compression ratio of 30:1 by wavelet transform; and (d) compression ratio of 30:1 by DCT-based JPEG algorithm.

514-61
 128447
 p. 4
 N93-19427

Maximal Codeword Lengths in Huffman Codes

Y. S. Abu-Mostafa

California Institute of Technology, Electrical Engineering Department

R. J. McEliece

Communications Systems Research Section

In this article, the authors consider the following question about Huffman coding, which is an important technique for compressing data from a discrete source. If p is the smallest source probability, how long, in terms of p , can the longest Huffman codeword be? It is shown that if p is in the range $0 < p \leq 1/2$, and if K is the unique index such that $1/F_{K+3} < p \leq 1/F_{K+2}$, where F_K denotes the K th Fibonacci number, then the longest Huffman codeword for a source whose least probability is p is at most K , and no better bound is possible. Asymptotically, this implies the surprising fact that for small values of p , a Huffman code's longest codeword can be as much as 44 percent larger than that of the corresponding Shannon code.

I. Introduction and Summary

Huffman coding is optimal (in the sense of minimizing average codeword length) for any discrete memoryless source, and Huffman codes are used widely in data compression applications. In many situations it would be useful to have an easy way to estimate the longest Huffman codeword length for a given source, without having to go through Huffman's algorithm, but since there is no known closed-form expression for the Huffman codeword lengths, no such estimate immediately suggests itself. However, since the longest codeword will always be associated with the least-probable source symbol, one way to address this problem is to ask the following question: If p is the smallest source probability, how long, in terms of p , can the longest Huffman codeword be? It turns out that this quantity, de-

noted by $L(p)$, is easy to calculate, and so $L(p)$ provides an "easy estimate" of the longest Huffman codeword length.

The formula for $L(p)$ involves the famous Fibonacci numbers $(F_n)_{n \geq 0}$, which are defined recursively, as follows:

$$F_0 = 0, F_1 = 1, \text{ and } F_n = F_{n-1} + F_{n-2} \text{ for } n \geq 2 \quad (1)$$

Thus, $F_2 = 1, F_3 = 2, F_4 = 3, F_5 = 5, F_6 = 8$, etc. The Fibonacci numbers and their properties are discussed in detail in [1, Section 1.2.8]. Here is the main result of this article. (Note that since the definition of $L(p)$ assumes p to be the smallest probability in a source, p must lie in the range $0 < p \leq 1/2$.)

Theorem 1. Let p be a probability in the range $0 < p \leq 1/2$, and let K be the unique index such that

$$\frac{1}{F_{K+3}} < p \leq \frac{1}{F_{K+2}} \quad (2)$$

Then $L(p) = K$. Thus $p \in (1/3, 1/2]$ implies $L(p) = 1$, $p \in (1/5, 1/3]$ implies $L(p) = 2$, $p \in (1/8, 1/5]$ implies $L(p) = 3$, etc.

It is easy to prove by induction that the Fibonacci numbers satisfy the following inequalities:

$$\phi^{n-2} < F_n < \phi^{n-1} \quad \text{for } n \geq 3 \quad (3)$$

where $\phi = (1 + \sqrt{5})/2 = 1.618\dots$ is the "golden ratio." By combining inequality (3) with Theorem 1, one sees that

$$\log_\phi \frac{1}{p} - 2 < L(p) < \log_\phi \frac{1}{p} \quad (4)$$

which, in turn, implies that

$$\lim_{p \rightarrow 0} \frac{L(p)}{\log_\phi \frac{1}{p}} = 1 \quad (5)$$

Since $\log_\phi x = (\log_2 x)/(\log_2 \phi) = 1.4404 \log_2 x$, Eq. (5) implies the surprising fact that for small values of p , a Huffman code's longest codeword can be as much as 44 percent larger than that of the corresponding (in general, suboptimal) Shannon code [2, Chapter 5], which assigns a symbol with probability p a codeword of length $\lceil \log_2 \frac{1}{p} \rceil$.

Theorem 1 is closely related to a result of Katona and Nemetz [4], which identifies the length of the longest possible Huffman codeword for a source symbol of probability p (whether or not p is the smallest source probability). Denoting this quantity by $L^*(p)$, their result is as follows:

Theorem 2. (Katona and Nemetz [4]) Let p be a probability in the range $0 < p < 1$, and let K be the unique index such that

$$\frac{1}{F_{K+2}} \leq p < \frac{1}{F_{K+1}} \quad (6)$$

Then $L^*(p) = K$. Thus, $p \in [1/2, 1)$ implies $L^*(p) = 1$, $p \in [1/3, 1/2)$ implies $L^*(p) = 2$, $p \in [1/5, 1/3)$ implies $L^*(p) = 3$, etc.

By comparing Theorems 1 and 2, one sees that $L^*(p) = L(p) + 1$ unless p is the reciprocal of a Fibonacci number, in which case $L^*(p) = L(p)$.¹

II. Proof of Theorem 1

The proof of Theorem 1 is in two parts. First, it will be shown that if $p > 1/F_{K+3}$, then in any Huffman code for a source whose smallest probability is p , the longest codeword length is at most K . In fact, a considerably stronger result will be proved. The class of efficient prefix codes will be defined, and it will be shown that any Huffman code, and in fact any optimal code for a given source, is efficient. Then it will be shown that if $p > 1/F_{K+3}$, in any efficient code for a source whose smallest probability is p , the longest codeword length is at most K . In the second half of the proof, it will be shown that if $p \leq 1/F_{K+2}$, there exists a source whose smallest probability is p , which has at least one Huffman code whose longest word has length K . As an extension, it will be seen that if $p < 1/F_{K+2}$, there exists a source whose smallest probability is p , and for which every optimal code has the longest word of length K . (If $p = 1/F_{K+2}$, however, there is no such source.)

Now comes the definition of efficient prefix codes, which is best stated in terms of the associated binary code tree (see Fig. 1). Each source symbol and its corresponding codeword is associated with a unique terminal node on the tree. Also, each node in the tree is assigned a probability. The probability of a terminal node is defined to be the probability of the corresponding source symbol, and the probability of any other node of the code tree is defined to be the sum of the probabilities of its two "children." The level of the root node is defined to be zero, and the level of every other node is defined to be one more than the level of its parent. Two nodes descended from the same parent node are called *siblings*. Figure 1 shows two different code trees for the source $[3/20, 3/20, 3/20, 3/20, 8/20]$. The tree in Fig. 1(a) corresponds to the prefix code $\{000, 001, 01, 10, 11\}$, and the tree in Fig. 1(b) corresponds to $\{000, 001, 010, 011, 1\}$.

Definition. A prefix code for a source S is efficient if every node except the root in the code tree has a sibling, and if $\text{level}(v) < \text{level}(v')$ implies $p(v) \geq p(v')$.

¹ In fact, however, if one were to make a subtle change in the definition of $L(p)$, this special case would disappear. The change required is to define $L(p)$ as the minimum maximum Huffman codeword length over all Huffman codes for a source with p as the least probability, where the outer minimum is over all Huffman codes for a given source.

Gallager [3] noted that every Huffman tree is efficient, but in fact it is easy to see more generally that every optimal tree is efficient. This is because in an *inefficient* tree, with nodes v and v' such that $\text{level}(v) < \text{level}(v')$ but $p(v) < p(v')$, by interchanging the subtrees rooted at v and v' , one arrives at a new code tree for the same source, whose average length has been reduced by exactly $(\text{level}(v') - \text{level}(v))(p(v') - p(v))$. However, it is not true that every efficient code is optimal. Indeed, Fig. 1 shows two different efficient code trees for the source $[3/20, 3/20, 3/20, 3/20, 8/20]$. The code in Fig. 1(b) is optimal, but the one in Fig. 1(a) is not.

Theorem 3. If $p > 1/F_{K+3}$, then in any efficient prefix code for a source whose least probability is p , the longest codeword length is at most K .

Proof: The contrapositive will be proved, i.e., if p is the least probability in a source that has an efficient prefix code whose longest word has length $\geq K + 1$, then $p \leq 1/F_{K+3}$.

Thus, suppose that S is a source whose least probability is p and that there is an efficient prefix code for S whose longest word is of length $\geq K + 1$. In the code tree for this code, there must be a path of length $K + 1$ starting from the terminal node, which corresponds to the longest word and moves upward toward the root. This path is shown in Fig. 2 as the path whose probabilities are p_0, p_1, \dots, p_{K+1} . Since the code is assumed to be efficient, each of the vertices in this path (except possibly the top vertex) has a sibling; these siblings are shown in Fig. 2 as having probabilities q_0, q_1, \dots, q_K . Now one can prove the following:

$$p_i \geq F_{i+2}p \quad \text{for } i = 0, 1, \dots, K + 1 \quad (7)$$

The proof of (7) is by induction. For $i = 0$, (7) merely says that $p_0 \geq p$, which is true since $p_0 = p$, by definition. Also, note that $q_0 \geq p$ since p is the least source probability. Thus, $p_1 = p_0 + q_0 \geq p + p = 2p = F_3p$, which proves (7) for $i = 1$. For $i \geq 2$, one has $p_i = p_{i-1} + q_{i-1}$. But $p_{i-1} \geq F_{i+1}p$ by induction, and $q_{i-1} \geq p_{i-2}$ since the code is efficient (q_{i-1} is a higher level node than p_{i-2}). Thus, one has $q_{i-1} \geq p_{i-2} \geq F_i p$ by induction, and so $p_i = p_{i-1} + q_{i-1} \geq (F_{i+1} + F_i)p = F_{i+2}p$, which completes the proof of (7).

Now consider the probability p_{K+1} . On one hand, $p_{K+1} \leq 1$; but on the other hand, $p_{K+1} \geq F_{K+3}p$, by (7). Thus, $p \leq 1/F_{K+3}$, which completes the proof. \square

Theorem 4. If $p \leq 1/F_{K+2}$, there exists a source whose smallest probability is p and which has a Huffman code whose longest word has length K . If $p < 1/F_{K+2}$, there exists such a source for which every optimal code has a longest word of length K .

Proof: Consider the following set of $K + 1$ source probabilities:

$$\left[p, \frac{F_1}{F_{K+2}}, \frac{F_2}{F_{K+2}}, \dots, \frac{F_{K-1}}{F_{K+2}}, \frac{F_K + 1}{F_{K+2}} - p \right] \quad (8)$$

Note that p is the minimal probability for this source, since $p \leq 1/F_{K+2} = F_1/F_{K+2}$. Now, consider the code tree for this source depicted in Fig. 3, which assigns the source probability p a word of length K . This tree is in fact a Huffman tree for these probabilities, i.e., a code tree that arises when Huffman's algorithm is applied to the source of (8). To see this, one first proves that the internal vertex probabilities p_i in Fig. 3 are given by the following formula:

$$p_i = F_{i+2}/F_{K+2} - h, \quad \text{for } i = 0, 1, \dots, K - 1 \quad (9)$$

$$p_K = 1 \quad (10)$$

where $h = 1/F_{K+2} - p$.

To prove (9), one uses induction. For $i = 0$, by definition, $p_0 = p = 1/F_{K+2} - h = F_2/F_{K+2} - h$. For $i \geq 1$, one then has $p_i = p_{i-1} + F_i/F_{K+2} = (F_{i+1}/F_{K+2} - h) + F_i/F_{K+2} = F_{i+2}/F_{K+2} - h$. To prove (10), note that $p_K = p_{K-1} + (F_K + 1)/F_{K+2} - p$. But from (9), $p_{K-1} = (F_K + 1)/F_{K+2} - h$, so that $p_K = (F_K + 1)/F_{K+2} - h + (F_K + 1)/F_{K+2} - p = F_{K+2}/F_{K+2} = 1$. Thus the probabilities in (8) sum to one.

It now follows that the tree in Fig. 3 is a Huffman tree, for from (9) one sees that at the i th stage ($i = 0, \dots, K - 1$), the "collapsed" source consists of the probabilities

$$\left[F_{i+2}/F_{K+2} - h, F_{i+1}/F_{K+2}, F_{i+2}/F_{K+2}, \dots, F_{K-1}/F_{K+2}, F_K/F_{K+2} + h \right] \quad (11)$$

Plainly the two leftmost probabilities in (11), namely $F_{i+2}/F_{K+2} - h$ and F_{i+1}/F_{K+2} , are two of the smallest probabilities, and so the tree of Fig. 3 is a Huffman tree, as asserted.

Finally, note that if $h > 0$, i.e., if $p < 1/F_{K+2}$, that the leftmost two probabilities in (11) are *uniquely* the two

smallest probabilities in the list, so that the Huffman tree in Fig. 3 is the unique Huffman tree for the source of Eq. (8). And since the set of codeword lengths in any optimal code is the same as the set of lengths in some Huffman code, the last statement in Theorem 4 follows. \square

By combining Theorems 3 and 4, one obtains a result that is stronger than Theorem 1.

Example 1: Let $p = 2^{-8}$. Then $1/F_{14} = 1/377 < p < 1/F_{13} = 1/233$, and so by Theorem 1, $L(2^{-8}) = 11$. More concretely, Theorem 3 shows that no Huffman code for a source whose smallest probability is 2^{-8} can have a codeword whose length is longer than 11. By Theorem 4, on the other hand, every optimal code for the source

$$\left[2^{-8}, \frac{1}{233}, \frac{1}{233}, \frac{2}{233}, \frac{3}{233}, \frac{5}{233}, \frac{8}{233}, \frac{13}{233}, \frac{21}{233}, \frac{34}{233}, \frac{55}{233}, \frac{90}{233} - 2^{-8} \right] \quad (12)$$

has a longest word of length 11. \square

III. Extension of the Katona-Nemetz Theorem

In this section, two theorems are stated without proof. When taken together, they yield a result that is slightly stronger than Katona and Nemetz's Theorem 2. The

proofs are entirely similar to the proofs of Theorems 3 and 4.

Theorem 5. Let S be a source containing a symbol a whose probability is p . If $p \geq 1/F_{K+2}$, then in any efficient prefix code for S , the length of the codeword assigned to the symbol a is at most K .

Theorem 6. Let $p < 1/F_{K+1}$. Then there exists a source S containing a symbol a whose probability is p , and such that every optimal code for S assigns a a codeword of length K . Explicitly, one such source is given by

$$S = \left[\frac{1}{F_{K+1}} - p - \epsilon, p, \frac{F_1}{F_{K+1}}, \frac{F_2}{F_{K+1}}, \dots, \frac{F_{K-1}}{F_{K+1}} + \epsilon \right] \quad (13)$$

where ϵ is any real number such that $0 < \epsilon < 1/F_{K+2} - p$.

Example 2: Let $p = 2^{-8}$. Then $1/F_{14} = 1/377 < p < 1/F_{13} = 1/233$, and so by Theorem 2, $L^*(2^{-8}) = 12$. Indeed, by Theorem 6, every optimal code for the source

$$\left[\frac{1}{233} - 2^{-8} - \epsilon, 2^{-8}, \frac{1}{233}, \frac{1}{233}, \frac{2}{233}, \frac{3}{233}, \frac{5}{233}, \frac{8}{233}, \frac{13}{233}, \frac{21}{233}, \frac{34}{233}, \frac{55}{233}, \frac{89}{233} + \epsilon \right] \quad (14)$$

where $0 < \epsilon < 1/233 - 1/256$, assigns the symbol with probability 2^{-8} a codeword of length 12. \square

Acknowledgment

The authors are grateful to Douglas Whiting of STAC, Inc., for suggesting this problem.

References

- [1] D. E. Knuth, *The Art of Computer Programming, vol. 1: Fundamental Algorithms*, 2nd ed., Reading, Massachusetts: Addison-Wesley, 1973.
- [2] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, New York: Wiley & Sons, 1991.
- [3] R. V. Gallager, "Variations on a theme by Huffman," *IEEE Trans. Inform. Theory*, vol. IT-24, pp. 668-674, November 1978.
- [4] G. O. H. Katona and T. O. H. Nemetz, "Huffman Codes and Self-Information," *IEEE Trans. Inform. Theory*, vol. IT-22, pp. 337-340, May 1976.

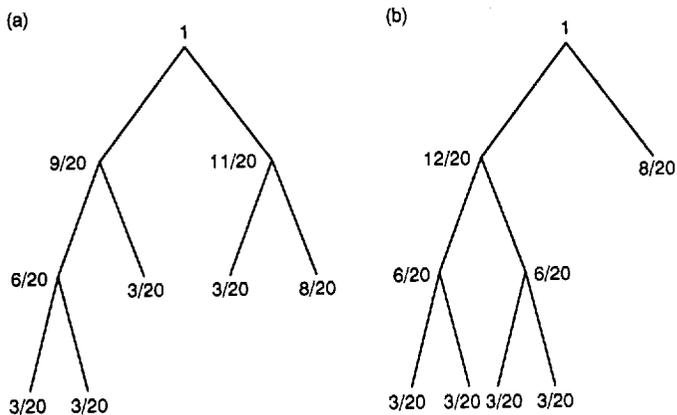


Fig. 1. Two code trees for the source $[3/20, 3/20, 3/20, 3/20, 3/20]$: (a) a tree that is efficient but not optimal (average length = 2.3) and (b) a tree that is optimal (average length = 2.2).

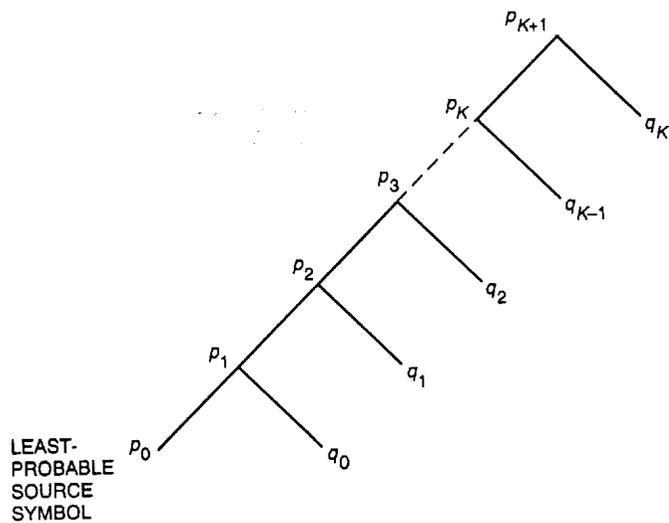


Fig. 2. A portion of an efficient code tree, in which the longest codeword has length $\geq K + 1$. p_0 is the least source probability.

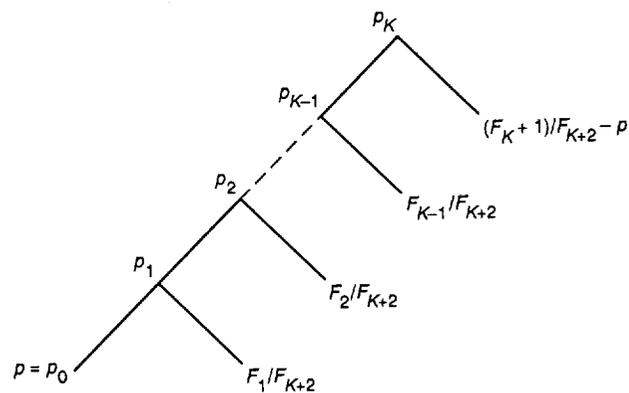


Fig. 3. A Huffman code tree for the source in (8). Its smallest probability is p , where $p \leq 1/F_{K+2}$, and its longest codeword length is K .

S/S - N'93 - 19428

128448

8
P

Comparisons of Theoretical Limits for Source Coding With Practical Compression Algorithms

F. Pollara and S. Dolinar

Communications Systems Research Section

In this article, the performance achieved by some specific data compression algorithms is compared with absolute limits prescribed by rate distortion theory for Gaussian sources under the mean square error distortion criterion. These results show the gains available from source coding and can be used as a reference for the evaluation of future compression schemes. Some current schemes perform well, but there is still room for improvement.

I. Introduction

The theoretical limits on the performance of source and channel coding are well known for several source and channel models [1,2,5]. In this article, the authors calculate the theoretical limits for one- and two-dimensional Gauss-Markov sources used as models for planetary images. The formulas underlying these calculations are well known; the aims in this article are first to collect and graphically display these results, and then to compare them with the performance of specific data compression algorithms.

These results show the gains available from source coding and can be used as a reference for the evaluation of present and future compression schemes. These results also suggest that large improvements in information transmission in future missions can be achieved by advanced source coding.

II. Theoretical Rate Distortion Limits

The authors consider time-discrete continuous-amplitude sources that produce identically distributed

output samples x governed by a probability distribution $P(x)$ with density $p(x)$. Each source sample x is reconstructed after source coding and decoding into a reconstructed sample y . The accuracy of reproduction is measured by a nonnegative function $d(x, y) = (x - y)^2$ called a squared error distortion measure. The average distortion D on a sequence of N samples is $(1/N) \sum_{i=0}^{N-1} (x_i - y_i)^2$ and is called mean square error (MSE) distortion.

A. One-Dimensional Gaussian Sources

For a Gaussian memoryless source, $p(x)$ is the Gaussian probability density with variance σ_x^2 , and the rate distortion function for MSE distortion is [1]

$$R(D) = \frac{1}{2} \log_2 \frac{\sigma_x^2}{D}, \quad 0 \leq D \leq \sigma_x^2 \quad (1)$$

where the rate R is measured in bits/sample.

A time-discrete stationary Gaussian source with spectral density function

$$\Phi(\omega) = \sum_{k=-\infty}^{\infty} \phi(n)e^{-jn\omega} \quad (2)$$

where $\phi(n)$ is the autocorrelation function, has a rate distortion $R(D)$ given in the parametric form [1]

$$D(\theta) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \min[\theta, \Phi(\omega)] d\omega \quad (3)$$

and

$$R(\theta) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \max \left[0, \frac{1}{2} \log_2 \frac{\Phi(\omega)}{\theta} \right] d\omega \quad (4)$$

where θ is the parameter.

Consider the special case of a first-order Gauss-Markov source of variance σ_x^2 with samples

$$x_i = \rho x_{i-1} + w_i \quad (5)$$

where $\{w_i\}$ is an independent, identically distributed (i.i.d.) zero-mean Gaussian sequence with variance $\sigma_w^2 = \sigma_x^2(1 - \rho^2)$. This source will be called the one-dimensional causal model, or 1DC model, and is characterized by an exponentially decaying memory given by the autocorrelation function

$$\phi(n) = \sigma_x^2 |\rho|^n, \quad 0 \leq \rho < 1 \quad (6)$$

which gives

$$\Phi(\omega) = \frac{\sigma_x^2(1 - \rho^2)}{1 - 2\rho \cos \omega + \rho^2} \quad (7)$$

Incidentally, the power spectral density function is always easy to find, given the definition of the model that generates the samples $\{x_i\}$, as described in [3].

B. Two-Dimensional Gaussian Sources

The rate distortion function $R(D)$ for a two-dimensional Gaussian source is given by [7]

$$D(\theta) = \left(\frac{1}{2\pi}\right)^2 \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} \min[\theta, \Phi(\omega_1, \omega_2)] d\omega_1 d\omega_2 \quad (8)$$

and

$$R(\theta) = \left(\frac{1}{2\pi}\right)^2 \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} \max \left[0, \frac{1}{2} \log_2 \frac{\Phi(\omega_1, \omega_2)}{\theta} \right] d\omega_1 d\omega_2 \quad (9)$$

A two-dimensional Gauss-Markov (autoregressive) causal source is defined by

$$x_{i,j} = \rho_1 x_{i-1,j} + \rho_2 x_{i,j-1} + \rho_{1,2} x_{i-1,j-1} + w_{i,j} \quad (10)$$

where $\{w_{i,j}\}$ is a two-dimensional i.i.d. zero-mean Gaussian sequence with variance σ_w^2 . If $\rho_{1,2} = -\rho_1 \rho_2$ is chosen, the source model in Eq. (10) becomes separable and will be called the two-dimensional causal (2DC) model. Then the variances of the sequences $\{w_{i,j}\}$ and $\{x_{i,j}\}$ are related by $\sigma_w^2 = \sigma_x^2(1 - \rho_1^2)(1 - \rho_2^2)$, and

$$\Phi(\omega_1, \omega_2) = \frac{\sigma_x^2(1 - \rho_1^2)(1 - \rho_2^2)}{(1 - 2\rho_1 \cos \omega_1 + \rho_1^2)(1 - 2\rho_2 \cos \omega_2 + \rho_2^2)} \quad (11)$$

This causal separable model has an autocorrelation function $\phi(n_1, n_2)$ given by

$$\phi(n_1, n_2) = \sigma_x^2 |\rho_1|^{n_1} |\rho_2|^{n_2} \quad (12)$$

which displays an undesirable nonisotropic behavior, as discussed later in this section.

Figure 1 shows the rate distortion functions for the 1DC model and the 2DC model with $\rho_1 = \rho_2 = \rho$ for several values of ρ . The values of the correlation coefficient ρ have been chosen to illustrate the effect of correlation on the rate necessary to represent the source. At low distortion, these values give rate distortion curves spaced by an integer number of bits from the curve for the memoryless source. Each successive correlation value in Fig. 1 represents (asymptotically for low distortion) one extra bit of information that can be extracted from each sample's correlation with its neighbors, and thus need not be spent to represent the source.

A more realistic model for images, the two-dimensional noncausal (2DNC) model, is given by

$$x_{i,j} = \alpha(x_{i,j-1} + x_{i,j+1} + x_{i-1,j} + x_{i+1,j}) + w_{i,j}, \quad |\alpha| < 1/4 \quad (13)$$

This is a noncausal model with a power spectral density

$$\Phi(\omega_1, \omega_2) = \frac{\sigma_w^2}{[1 - 2\alpha(\cos \omega_1 + \cos \omega_2)]^2} \quad (14)$$

where $\sigma_w^2 = \sigma_x^2 \eta$, and

$$\eta^{-1} = \left(\frac{1}{2\pi}\right)^2 \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} \frac{1}{[1 - 2\alpha(\cos \omega_1 + \cos \omega_2)]^2} d\omega_1 d\omega_2 \quad (15)$$

The autocorrelation of this model was computed numerically, and it was found to behave almost isotropically at small displacements. Figure 2 shows a comparison of the autocorrelations for the causal and noncausal models for correlation coefficients ρ that are representative of typical planetary images. The function

$$\phi(n_1, n_2) = \sigma_x^2 |\rho| \sqrt{n_1^2 + n_2^2} \quad (16)$$

is an example of exactly isotropic autocorrelation [4], but the authors do not presently know a model that realizes such an autocorrelation.

The qualitative behavior of the autocorrelation functions for the 2DC and 2DNC models is illustrated in the contour plots of Fig. 3. Note that for small values of n_1 and n_2 the contours for the 2DNC model are nearly circular, indicating that this model is nearly isotropic for small displacements.

The rate distortion functions for the 2DC and 2DNC models are shown in Fig. 4 with the same parameter values used in Fig. 2. Since the autocorrelation function for the 2DNC model decays more rapidly than for the 2DC model when both models have the same value of $\phi(1, 0)$ fixed, the rate distortion function of the 2DNC model lies above that of the 2DC model.

III. Quantization

Given a time-discrete continuous amplitude source, the simplest form of data compression is scalar (sample-by-sample) quantization. An M -level quantizer is a device with an input that can assume any real value x and an output y that can assume only M values $\{L_1, \dots, L_M\}$. Usually, the number of levels is a power of 2, so that a B -bit quantizer has 2^B levels. Given the quantization thresholds

$\{T_1, \dots, T_{M-1}\}$, the output is $y = L_k$ if and only if $T_{k-1} < x \leq T_k$, $k = 1, \dots, M$, where $T_0 = -\infty$ and $T_M = +\infty$. The input-output characteristic of a four-level quantizer is shown in Fig. 5.

Let $\{x_1, \dots, x_N\}$ be a sequence of random samples generated by a source and let $\{y_1, \dots, y_N\}$ be the corresponding quantized samples produced by an M -level quantizer. Then the quantized sequence has rate $B = \log_2 M$ bits and MSE distortion

$$\begin{aligned} D &\triangleq \frac{1}{N} \sum_{i=1}^N E[(x_i - y_i)^2] = E[(x_i - y_i)^2] \\ &= \sum_{k=1}^M \int_{T_{k-1}}^{T_k} (x - L_k)^2 p(x) dx \end{aligned} \quad (17)$$

where $p(x)$ is the probability density of the source. Therefore, the M -level quantizer realizes the point (B, D) on the rate distortion plane. The optimum quantizer, which achieves the lowest possible MSE for given source statistics, has been determined in terms of the reproduction levels $\{L_k\}$ and the thresholds $\{T_k\}$ using an optimization technique developed by Lloyd and Max in 1960. If the quantizer is restricted to have equally spaced thresholds, i.e., a uniform quantizer with constant step size $T_k - T_{k-1}$ is considered, a slightly higher distortion for corresponding rates is obtained, as shown in Fig. 6 for the Gaussian memoryless source. An optimum uniform quantizer is a uniform quantizer that minimizes the MSE distortion.

Improved rate performance can be obtained by using entropy coding after quantization, since the probability $P_k = \Pr(y = L_k) = \int_{T_{k-1}}^{T_k} p(x) dx$ that a quantizer output will be L_k is not a constant (except for degenerate cases), and therefore the entropy of the quantized samples y is strictly less than B

$$H(y) = - \sum_{k=1}^N P_k \log_2 P_k < B \quad (18)$$

The entropy coded performance of the two quantizers considered above is also shown in Fig. 6, where it is apparent that the advantage of the Lloyd-Max quantizer over the uniform quantizer disappears after entropy coding. Results on entropy coded quantizers were obtained from the literature [3] and reproduced by computer simulation.

Instead of quantizing individual source samples, one could collect a whole vector $\mathbf{x} = (x_1, \dots, x_n)$ and then

vector quantization. The performance of vector quantization methods will be discussed in a future article.

If one replaces the memoryless Gaussian source with a one-dimensional Gauss-Markov source with correlation coefficient ρ between successive samples (1DC model), a simple method to exploit the source memory is to take differences between successive quantized samples and then apply entropy coding. The performance of such a one-step predictor on samples produced by an optimum uniform quantizer is shown in Fig. 7.

In practice, the continuous amplitude source is initially quantized to B bits, typically 8 bits. In the following discussion of practical compression algorithms for images, it is assumed that the source has been quantized to 8 bits per sample by an optimum uniform quantizer.

IV. Comparisons of Practical Compression Algorithms and Theoretical Limits

The performance of specific compression algorithms designed for 8-bit input data can be measured experimentally by generating in software a Gauss-Markov random field according to one of the models described in Section II and by quantizing the resulting samples to 8 bits with an optimum uniform quantizer.

The entropy coded one-step predictor described in the previous section is a simple example of a practical compression scheme, and it is essentially the image compression scheme used in Voyager, where the source was initially quantized to 8 bits by the camera. The point denoted by 8 in the rate distortion plot of Fig. 7 represents the so-called lossless performance of such a scheme. This scheme performs reasonably well at low distortions (as compared with the rate distortion function) when it is applied to the one-dimensional source 1DC. One will see that its performance is no longer attractive when applied to two-dimensional sources 2DC or 2DNC.

The proposed Joint Photographic Expert Group (JPEG) image compression standard [6] uses, in its baseline version, discrete cosine transform (DCT) processing, quantization, and Huffman coding. The performance of this compression scheme on the 2DC model has been evaluated and compared with the rate distortion limits in Fig. 8. The performance of the entropy coded one-step predictor on the 2DC model is also shown in Fig. 8 for comparison. For most science purposes, a typical planetary image is considered acceptable at normalized distortions D/σ_x^2 up to approximately 10^{-2} , corresponding to about 5 gray levels of rms error out of 256 levels for typical images. In this range of interest, the JPEG scheme is superior to the entropy coded predictor, but the theoretical limit leaves ample space for improvements. The performances of the JPEG scheme and the entropy coded one-step predictor on the 2DNC model are compared in Fig. 9.

V. Conclusion

The theoretical limits computed in this article and the experimental results on source models verify the gains available by source coding, and can be used as a reference for the evaluation of present and future compression schemes. These results also suggest that large improvements in information transmission in future missions can be achieved by advanced source coding.

Mathematical source models studied in this article include both relatively simple one- and two-dimensional causal Gauss-Markov models and a two-dimensional non-causal model whose nearly isotropic correlation function more closely resembles that of real images.

More work is necessary in relating the mathematical models to actual image sources, in evaluating the performance of other practical compression schemes, and in understanding the actual quantization performed in the camera.

References

- [1] T. Berger, *Rate Distortion Theory*, Englewood Cliffs, New Jersey: Prentice Hall, 1971.
- [2] S. Dolinar and F. Pollara, "The Theoretical Limits of Source and Channel Coding," *TDA Progress Report 42-102*, vol. April-June 1990, Jet Propulsion Laboratory, Pasadena, California, pp. 62-72, August 15, 1990.
- [3] A. K. Jain, *Fundamentals of Digital Image Processing*, Englewood Cliffs, New Jersey: Prentice Hall, 1989.
- [4] N. S. Jayant and P. Noll, *Digital Coding of Waveforms*, Englewood Cliffs, New Jersey: Prentice Hall, 1984.
- [5] R. J. McEliece, *The Theory of Information and Coding*, Reading, Massachusetts: Addison Wesley, 1977.
- [6] F. Pollara and S. Arnold, "Emerging Standards for Still Image Compression: A Software Implementation and Simulation Study," *TDA Progress Report 42-104*, vol. October-December 1990, Jet Propulsion Laboratory, Pasadena, California, pp. 98-102, February 15, 1991.
- [7] J. A. Stuller and B. Kurz, "Intraframe Sequential Picture Coding," *IEEE Transactions on Communications*, vol. COM-25, no. 5, pp. 485-495, May 1977.

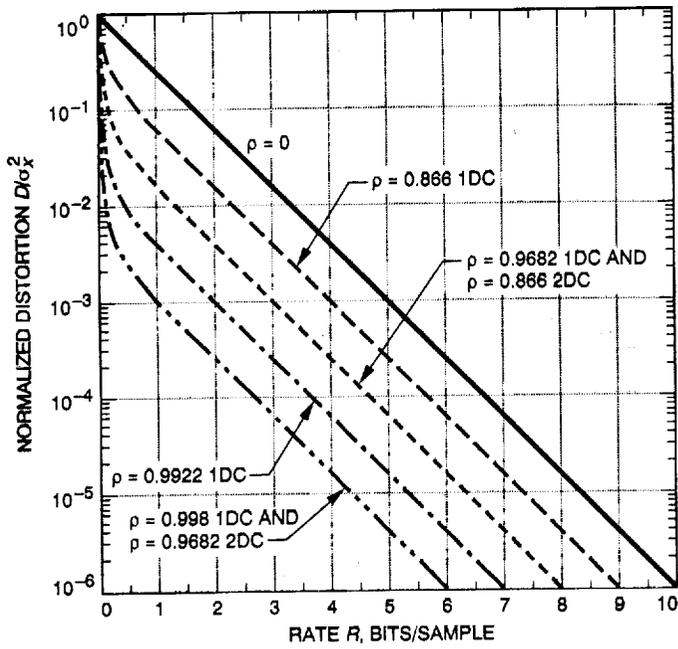


Fig. 1. Rate distortion functions for 1DC and 2DC models.

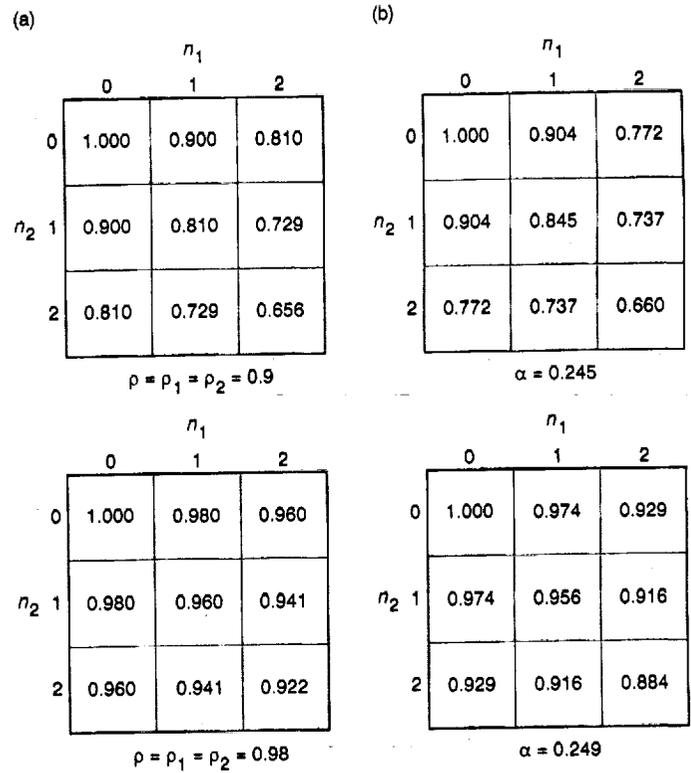


Fig. 2. Two-dimensional normalized autocorrelation functions $\phi(n_1, n_2)/\sigma_x^2$: (a) 2DC model and (b) 2DNC model.

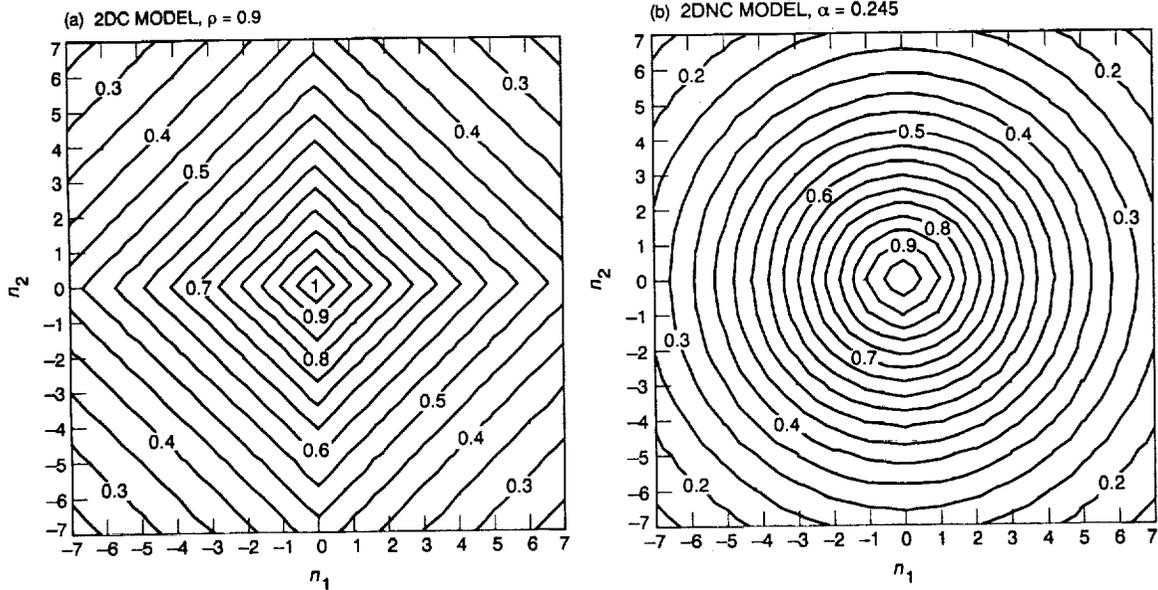


Fig. 3. Contour plots: (a) causal and (b) noncausal autocorrelations.

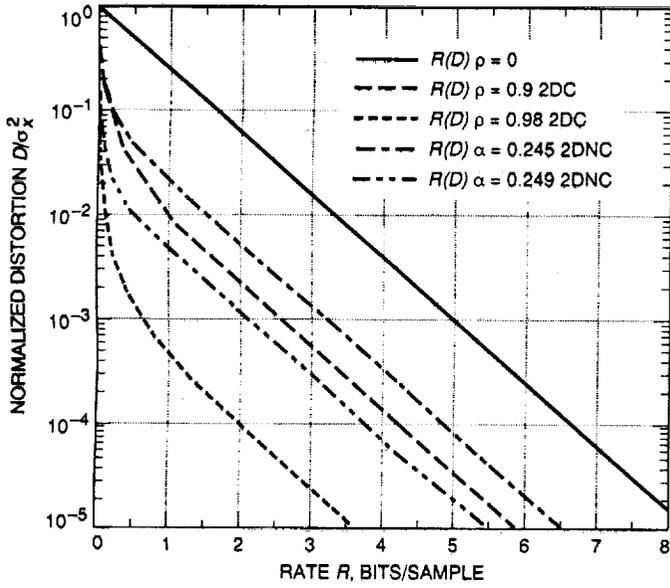


Fig. 4. Rate distortion functions for causal and noncausal two-dimensional models.

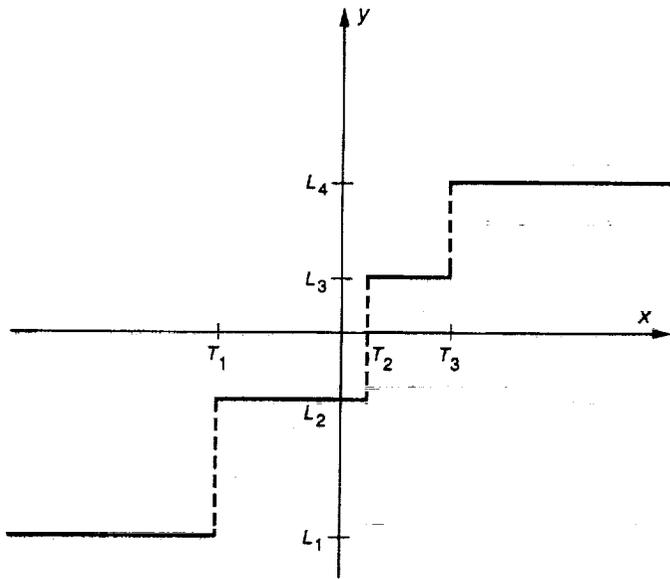


Fig. 5. Input-output characteristic of a four-level quantizer.

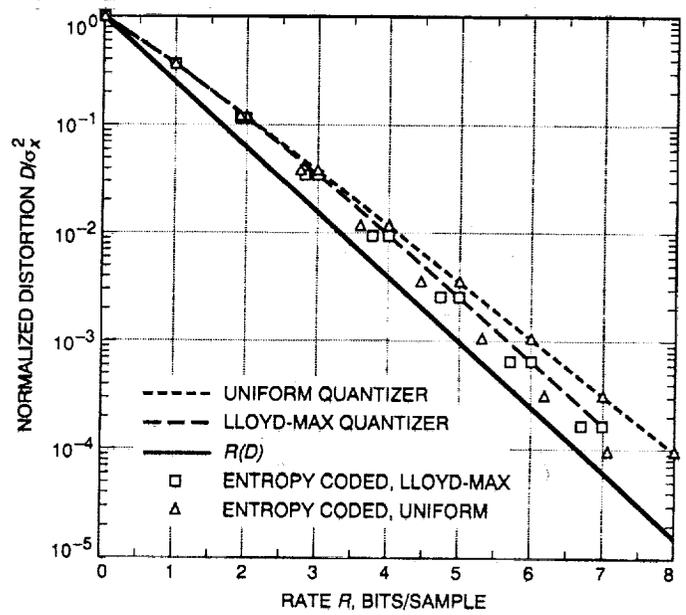


Fig. 6. Performance of quantization schemes for memoryless Gaussian sources.

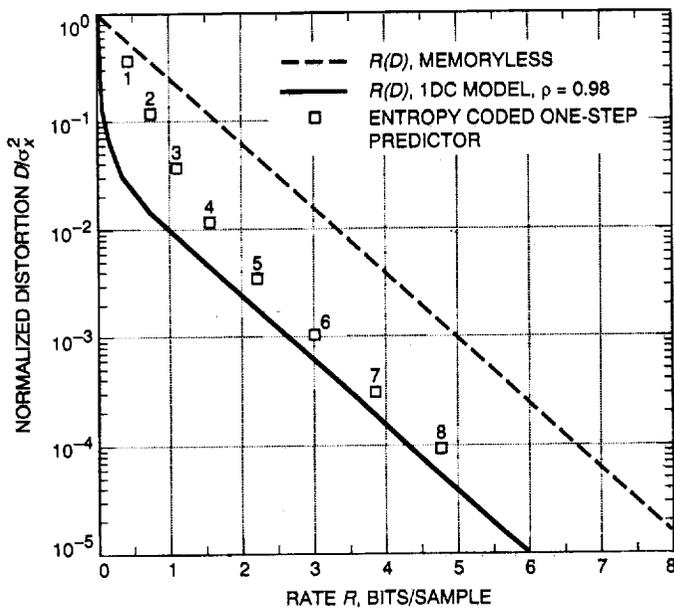


Fig. 7. Performance of entropy coded one-step predictor on Gauss-Markov source with $\rho = 0.98$.

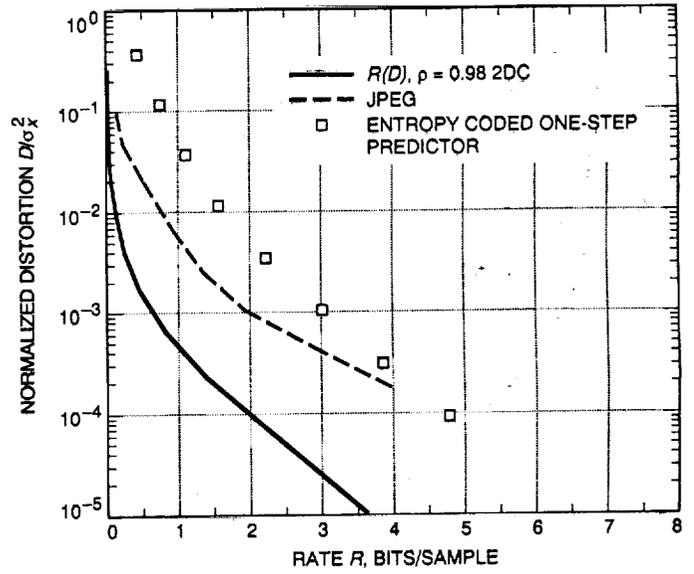


Fig. 8. Comparisons of practical compression algorithms and theoretical limits (2DC model).

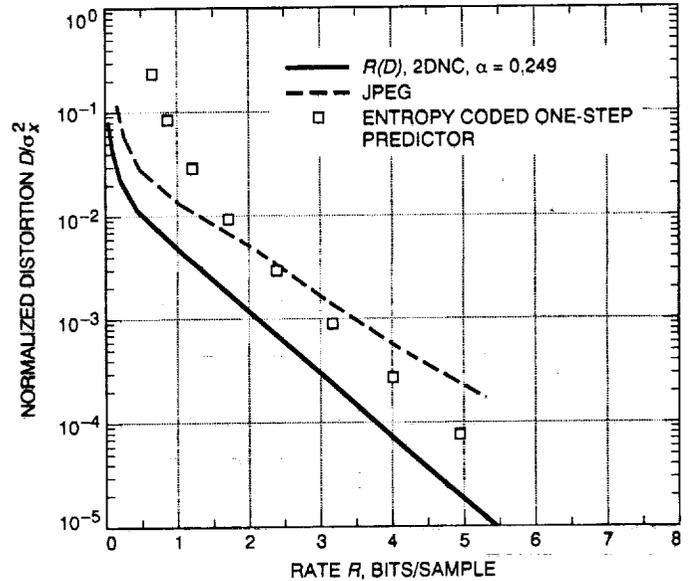


Fig. 9. Comparisons of practical compression algorithms and theoretical limits (2DNC model).

516N93-19429

128449

P=6

Cascaded Convolutional Codes

F. Pollara and D. Divsalar
Communications Systems Research Section

Due to the hardware design of Galileo's Command and Data Subsystem (CDS), the channel code usable in an S-band (2290-2300 MHz) mission must include the NASA standard (7,1/2) convolutional code. Galileo's hardware encoder for the (15,1/4) code is not usable in S-band mode. However, the need for higher coding gain dictates the use of long constraint length convolutional codes. Theoretical results show how a large subclass of such codes is realizable by using a software encoder in the CDS cascaded with the hardware encoder for the NASA standard code.

I. Introduction

Several options for improving Galileo's telemetry downlink performance at Jupiter if the high-gain antenna fails to deploy were evaluated in the *Galileo Options Study*¹ sponsored by the Telecommunications and Data Acquisition (TDA) Office. Specific recommendations were developed in the subsequent *Galileo S-Band Mission Study*.²

In this article, the authors describe one of the proposed options to improve Galileo's S-band (2290-2300 MHz) downlink performance based on the use of advanced long constraint length convolutional codes.

The Command and Data Subsystem (CDS) of Galileo provides two output paths to the Modulation/Demodulation Subsystem (MDS): a low-rate telemetry output

(40 bps) and a high-rate telemetry output (10 bps to 134.4 kbps). The low-rate output is directly connected to the low-gain antenna path. The high-rate output may use the low-gain antenna only through a hardware (7,1/2) convolutional encoder, as shown in Fig. 1. Galileo's hardware encoder for the (15,1/4) code is not usable in S-band mode.

One of the options to improve Galileo data return through its low-gain antenna is to use advanced channel coding techniques, including long constraint length convolutional codes. This could be achieved by uploading a software (15,1/4) encoder in the CDS and using the low-rate output connected to the S-band telemetry path. However, the encoder output rate would be fixed to 40 symbols per second, which is not compatible with the desire for higher rates.

The only way to send S-band telemetry at higher rates is to use the high-rate CDS output, which then forces the use of the hardware (7,1/2) convolutional encoder. Therefore, methods are to be sought for realizing long constraint

¹L. Deutsch, *Galileo Options Study* (internal document), Jet Propulsion Laboratory, Pasadena, California, November 5, 1991.

²L. Deutsch and J. Marr, *Galileo S-Band Mission Study Final Report* (internal document), Jet Propulsion Laboratory, Pasadena, California, March 2, 1992.

length convolutional codes by cascading a software encoder with the existing hardware (7,1/2) encoder.

II. Cascaded Convolutional Codes

The best solution would be to find a method for bypassing the (7,1/2) hardware encoder by realizing an inverse software encoder preceding it. Any desired code would then be realizable in software.

It is well known that a noncatastrophic encoder has a feed-forward inverse [1]. Denote N as the (7,1/2) hardware encoder on Galileo with generator polynomials g_0 and g_1 , and N^{-1} as its inverse. Then it is possible to undo the operation of N , as shown in Fig. 2, where M is a multiplexer and the relative symbol rates are shown below each connection.

For the (7,1/2) standard NASA code having generator polynomials $g_0 = 1 + x + x^2 + x^3 + x^6$ and $g_1 = 1 + x^2 + x^3 + x^5 + x^6$, the feed-forward inverse is delay-free, since the greatest common divisor (GCD) $(g_0, g_1) = 1$, and is given by $f_0 = 1 + x + x^2 + x^3 + x^4$ and $f_1 = x^2 + x^4$, since $\text{GCD}(g_0, g_1) = g_0 f_0 + g_1 f_1 = 1$. The feed-forward inverse can be used to recover the information sequence a from the encoder output y . This inverse cannot be used in the Galileo CDS configuration since the output of the hardware encoder is not accessible.

However, the problem of constructing a preinverse of N such that the sequences w and y are identical, as shown in Fig. 3, has no solution in general. This is clear from an information-theory point of view, since w can be any binary sequence while the sequence y is restricted to being a code word of the specific code in use.

A. Structure and Design of Cascaded Codes

An alternative method for realizing a longer constraint length code with some form of processing preceding the hardware encoder is shown in Fig. 4.

The structure shown in Fig. 4 is just an example of one of the possible structures for obtaining a rate 1/4 code equivalent to code D , shown in Fig. 5. This structure cannot obtain all desired codes but just a certain subclass of codes.

A simplified strategy for designing a (15,1/4) cascaded code is to assume that a code C specified through f_0 and f_1 is given and then compute the resulting equivalent code specified by h_0, h_1, h_2 and h_3 . One has

$$y_0(x) = a(x)f_0(x)$$

$$y_1(x) = a(x)f_1(x)$$

$$z(x) = y_0(x^2) + xy_1(x^2)$$

$$u_0(x) = z(x)g_0(x)$$

$$u_1(x) = z(x)g_1(x)$$

$$w(x) = u_0(x^2) + xu_1(x^2)$$

which gives

$$w(x) = a(x^4)[f_0(x^4) + x^2 f_1(x^4)][g_0(x^2) + xg_1(x^2)]$$

and, from Fig. 5,

$$q(x) = s_0(x^4) + xs_1(x^4) + x^2 s_2(x^4) + x^3 s_3(x^4)$$

$$s_0(x) = a(x)h_0(x)$$

$$s_1(x) = a(x)h_1(x)$$

$$s_2(x) = a(x)h_2(x)$$

$$s_3(x) = a(x)h_3(x)$$

which gives

$$q(x) = a(x^4)[h_0(x^4) + xh_1(x^4) + x^2 h_2(x^4) + x^3 h_3(x^4)]$$

Since one wants

$$q(x) = w(x)$$

to hold, the condition becomes

$$\begin{aligned} & [f_0(x^4) + x^2 f_1(x^4)][g_0(x^2) + xg_1(x^2)] = \\ & [h_0(x^4) + xh_1(x^4) + x^2 h_2(x^4) + x^3 h_3(x^4)] \end{aligned} \quad (1)$$

By expanding the left-hand side of Eq. (1) and identifying terms of equal power, one can find polynomials h_0 , h_1 , h_2 , and h_3 satisfying this equation. The memory of a convolutional code is the maximum degree among its generator polynomials. If one defines by m_C , m_N , and m_D the memories of codes C , N , and D , Eq. (1) implies that

$$2m_D = 2m_C + m_N$$

The same result applies to the respective constraint lengths, since for these codes the constraint length K is equal to $m + 1$.

B. Example

In this example, it is assumed that g_0 , g_1 , f_0 , and f_1 are given and a solution for h_0 , h_1 , h_2 , and h_3 is sought. In particular, if one chooses

$$f_0(x) = 1 + x^9 + x^{11}$$

$$f_1(x) = 1 + x^3 + x^6 + x^9 + x^{11}$$

one gets $h_0 = 57347$, $h_1 = 71526$, $h_2 = 02245$, and $h_3 = 52207$, in octal representation. This yields a rate $r = 1/4$ code with memory $m = 14$ (constraint length $K = m + 1 = 15$), i.e., a (15,1/4) code with free distance $d_f = 30$, while the experimental code on Galileo had $d_f = 35$ with the same parameters. Searching over f_0 and f_1 may yield better codes.

A more explicit solution for Eq. (1) can be found by defining g_0 and g_1 in terms of their even and odd parts

$$\left. \begin{aligned} g_0(x) &\triangleq g_{0e}(x^2) + xg_{0o}(x^2) \\ g_1(x) &\triangleq g_{1e}(x^2) + xg_{1o}(x^2) \end{aligned} \right\} (2)$$

Then it follows that

$$\left. \begin{aligned} h_0(x) &= f_0(x)g_{0e}(x) + xf_1(x)g_{0o}(x) \\ h_1(x) &= f_0(x)g_{1e}(x) + xf_1(x)g_{1o}(x) \\ h_2(x) &= f_0(x)g_{0o}(x) + f_1(x)g_{0e}(x) \\ h_3(x) &= f_0(x)g_{1o}(x) + f_1(x)g_{1e}(x) \end{aligned} \right\} (3)$$

As a verification, the results obtained in the previous example can be reproduced by using this explicit solution.

Convolutional codes with high coding gain, including the original (15,1/4) Galileo code, are such that the first and last coefficient of all generator polynomials are equal to 1, i.e., $h_{i,j} = 1$, $i = 0, 1, 2, 3$, $j = 0, 14$, where $h_i(x) \triangleq \sum_{j=0}^{14} h_{i,j}x^j$. Therefore, it is interesting to determine whether a cascaded code having this property exists. From Eqs. (3) one has

$$f_{0,0}g_{0e,0} = h_{0,0}$$

$$f_{0,0}g_{1e,0} = h_{1,0}$$

$$f_{0,0}g_{0o,0} + f_{1,0}g_{0e,0} = h_{2,0}$$

$$f_{0,0}g_{1o,0} + f_{1,0}g_{1e,0} = h_{3,0}$$

where $f_i(x) \triangleq \sum_{j=0}^{11} f_{i,j}x^j$ and $g_{ip}(x) \triangleq \sum_{j=0}^6 g_{ip,j}x^j$. In order to get $h_{i,0} = 1$ and $i = 0, 1, 2, 3$, one should have

$$g_{0e,0} = g_{1e,0} = 1 \text{ and } g_{0o,0} = g_{1o,0}$$

Also, from Eq. (3) one has the following conditions, on the coefficients of x^{14}

$$f_{0,11}g_{0e,3} + f_{1,11}g_{0o,2} = h_{0,14}$$

$$f_{0,11}g_{1e,3} + f_{1,11}g_{1o,2} = h_{1,14}$$

$$f_{1,11}g_{0e,3} = h_{2,14}$$

$$f_{1,11}g_{1e,3} = h_{3,14}$$

In order to get $h_{i,14} = 1$ and $i = 0, 1, 2, 3$, one should have

$$g_{0e,3} = g_{1e,3} = 1 \text{ and } g_{0o,2} = g_{1o,2}$$

But, for the (7,1/2) NASA code, one has

$$g_{0e,0} = g_{1e,0} = 1$$

$$g_{0o,0} = 1 \text{ and } g_{1o,0} = 0$$

$$g_{0e,3} = g_{1e,3} = 1$$

$$g_{0o,2} = 0 \text{ and } g_{1o,2} = 1$$

$$h_{0,14} \neq h_{1,14} \text{ and } h_{2,14} = h_{3,14} = 1$$

which implies $g_{0o,0} \neq g_{1o,0}$ and $g_{0o,2} \neq g_{1o,2}$. Thus, it is impossible to get $h_{i,0} = 1$ for all i 's and/or $h_{i,14} = 1$ for all i 's. When the NASA code is used, the following is possible

$$h_{0,0} = h_{1,0} = 1 \text{ and } h_{2,0} \neq h_{3,0}$$

and

III. Conclusion

A method is presented for realizing long constraint length convolutional codes as a cascade of two codes including the NASA standard (7,1/2) code. This analysis shows that a large class of codes can be realized using this construction method. These results led to the inclusion of one of these cascaded codes in the design described in the *Galileo S-Band Mission Study*.

Reference

- [1] J. L. Massey and M. K. Sain, "Inverses of Linear Sequential Circuits," *IEEE Transactions on Computers*, vol. C-17, pp. 330-337, April 1968.

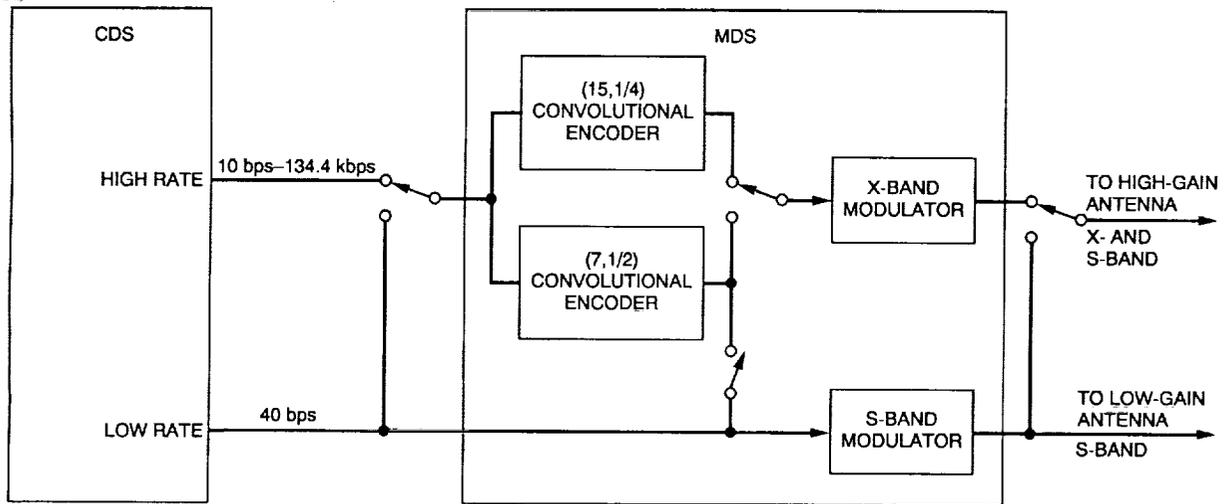


Fig. 1. Functional telemetry data flow.

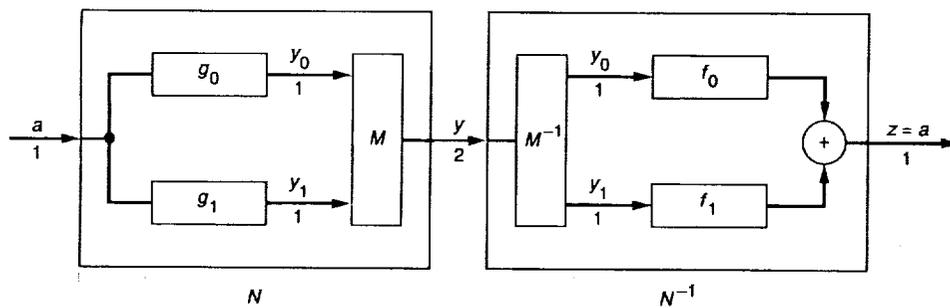


Fig. 2. Feed-forward inverse of code N .

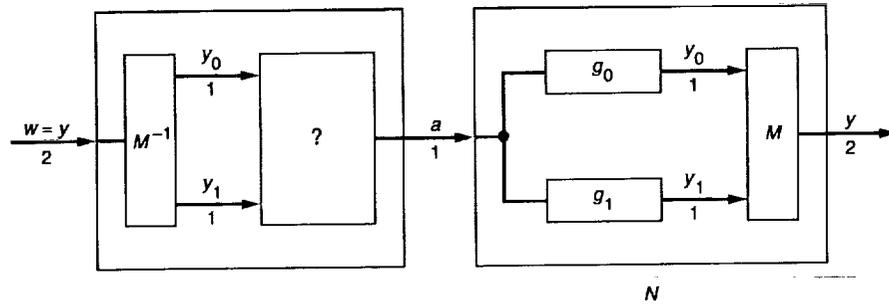


Fig. 3. Preinverse of code N .

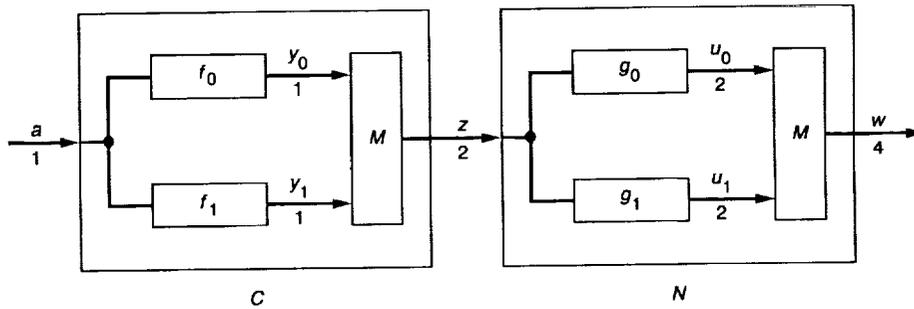


Fig. 4. Alternate structure for proposed Galileo code (rate = $1/4$).

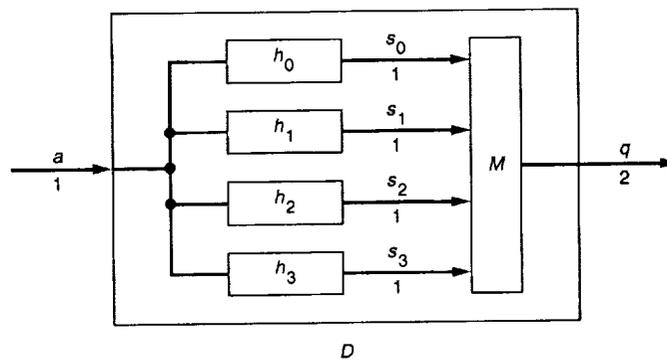


Fig. 5. Code equivalent to the cascades of codes C and N .

517-663-19430
 128450
 p. 8

Binary Weight Distributions of Some Reed-Solomon Codes

F. Pollara and S. Arnold

Communications Systems Research Section

The binary weight distributions of the (7,5) and (15,9) Reed-Solomon (RS) codes and their duals are computed using the MacWilliams identities. Several mappings of symbols to bits are considered and those offering the largest binary minimum distance are found. These results are then used to compute bounds on the soft-decoding performance of these codes in the presence of additive Gaussian noise. These bounds are useful for finding large binary block codes with good performance and for verifying the performance obtained by specific soft-decoding algorithms presently under development.

I. Introduction

Reed-Solomon (RS) codes are currently used in the DSN as outer codes in a concatenated coding system. For this application, they are decoded by algebraic techniques using operations in the field over which the code is designed. An (n, k) RS code C over $GF(2^m)$ has codewords of length $n = 2^m - 1$ symbols, where each symbol is a binary m -tuple. Let A_i be the number of codewords of weight i in C , then the vector (A_0, A_1, \dots, A_n) is called the weight distribution of C , where the weight (Hamming weight) of a codeword is the number of its nonzero coordinates. The term "coordinate" assumes different meanings depending on how one views the code: One may assume that there are n coordinates, each having a value in $GF(2^m)$, or one may consider the binary expansion of the code, i.e., a binary (nm, km) code, where each coordinate is a single bit. Hence, one may be interested in the symbol weight distribution or in the binary weight distribution of a (nonbinary) code. The latter depends on the specific symbol to binary m -tuple mapping that was chosen. Which of these distributions is of interest depends on which type

of decoding algorithm one plans to use, since weight distributions are essential in evaluating the error-correcting performance of a code. The symbol weight distribution of RS codes is well known [1] and can be used to find the performance of algebraic decoders working on symbols. The full error-correcting power of a code is obtained when soft, maximum-likelihood decoding is used, working directly on unquantized vectors in the nm -dimensional Euclidean space. Soft, maximum-likelihood decoding is superior to its hard quantized version by more than 2 dB. Furthermore, the algebraic decoding techniques usually employed for RS codes are not maximum-likelihood, but rather "incomplete" decoding techniques with a nonzero probability of decoding failure.

II. Binary Weight Distribution

This article focuses on evaluating the soft, maximum-likelihood decoding performance of RS codes, and therefore one needs to compute the binary weight enumerators of these codes. Such a task is a long-standing open prob-

lem in coding theory due to its intrinsic complexity. However, approximate results have been found and results for special classes of codes are known.

In general, one could think of using an exhaustive enumeration to find the numbers A_i by considering each codeword. Unfortunately, such a method is limited to fairly short codes, even on the most powerful computers available.

It was possible, for example, to find by exhaustive enumeration the weight distribution of a (21,15) binary code obtained from the (7,5) RS code over $\text{GF}(2^3)$, but it was impractical to find that of a (60,36) binary code obtained from the (15,9) RS code over $\text{GF}(2^4)$, since it involves 2^{36} codewords. Fortunately, a well-known result from coding theory, the MacWilliams identities [2], can be used to relate the weight distribution of a code to that of its dual. For example, one can find the binary weight distribution of the (15,9) RS code from that of its (15,6) dual code, by exhaustive enumeration on 2^{24} codewords instead of 2^{36} codewords.

Let the weight enumerator of a code C be defined as $W_C(x, y) = \sum_{i=0}^n A_i x^{n-i} y^i$. Then the weight enumerator of the dual code C^\perp of a binary code C is given by [MacWilliams identity over $\text{GF}(2)$]

$$W_{C^\perp} = \frac{1}{2^k} W_C(x + y, x - y)$$

The generator polynomial of an (n, k) RS code C may be written as

$$g(x) = \prod_{i=1}^{n-k} (x - \alpha^{i+b})$$

where b can be chosen among the values $0, 1, \dots, n-1$, and α is a root of the primitive polynomial over $\text{GF}(2)$ defining the field $\text{GF}(2^m)$. The parity check polynomial $h(x)$ of the code C

$$h(x) = \frac{x^n - 1}{g(x)} = \prod_{i=n-k+1}^n (x - \alpha^{i+b})$$

is the generator of the dual code C^\perp .

The binary weight distribution of the (21,15) binary code derived from the (7,5) RS code is shown in Table 1

together with the distribution of the (21,6) dual code associated with the (7,2) RS code. Results are shown for different values of the parameter b that correspond to different assignments of symbols to binary m -tuples. These are only a small subset of all possible assignments. The weight distributions shown in Table 1 could be found by exhaustive enumeration. For the (7,2) RS code, the largest binary minimum distance found was 8, which is the best possible according to [4]. For the (7,5) RS code the best result was $d_{min} = 4$, which meets the Griesmer upper bound [3].

The weight distribution of the (60,36) binary code was found by using the MacWilliams identity for binary codes, by a procedure shown in Fig. 1. First, the (15,6) dual code was generated by using the parity check polynomial of the (15,9) code as its generator. Then, the (15,6) code over $\text{GF}(2^4)$ was represented as a binary (60,24) code by mapping symbols in $\text{GF}(2^4)$ to binary 4-tuples by using the representation of field elements given by the irreducible polynomial $1 + x + x^4$ over $\text{GF}(2)$. The weight distribution of the (60,24) code was found by exhaustive enumeration, and finally, the weight distribution of the (60,36) code was computed by the MacWilliams identity for binary codes.

The missing arrow in the block diagram of Fig. 1 stresses the fact that the resulting (60,36) code is not necessarily related to its nonbinary parent, the (15,9) code, by the same mapping relating the (15,6) code to the (60,24) code. Table 2 shows the binary weight distributions for some (60,24) codes derived from the (15,6) RS code, where the largest minimum distance found was 13. It is known [4] that at least one (60,24) code exists for some value of d_{min} in the range 16 to 18. Table 3 shows similar results for the (60,36) code, where the largest minimum distance found was 8. At least one (60,36) code exists for some value of d_{min} in the range 9 to 12 [4].

III. Performance Evaluation

The soft decoding performance of block codes can be estimated by union bounding techniques. Specifically the word error probability P_w is upper bounded by [5]

$$P_w \leq \frac{1}{2} \sum_{j=2}^M \text{erfc} \left(\sqrt{w_j R \frac{E_b}{N_o}} \right)$$

where $R = k/n$ is the code rate, $M = 2^k$ is the number of codewords, and w_j is the weight of the j th codeword. The bound on P_w may be easily rewritten in terms of the weight distribution A_i as

$$P_w \leq \frac{1}{2} \sum_{i=1}^n A_i \operatorname{erfc} \left(\sqrt{iR \frac{E_b}{N_o}} \right)$$

Similarly, for hard quantized, maximum-likelihood decoding one can derive the union bound [5]

$$P_w \leq \sum_{j=2}^M \left[\sqrt{4p(1-p)} \right]^{w_j}$$

where $p = \frac{1}{2} \operatorname{erfc} \left(\sqrt{R \frac{E_b}{N_o}} \right)$.

The word error probability P_w can be related to the average bit error probability P_b by observing that when at least $t+1$ errors occur, the decoder produces an erroneous codeword containing at least $d_{min} = 2t + 1$ errors over n symbols. Therefore, kd_{min}/n is the average number of erroneous bits. Since in a codeword there are k bits, one has

$$P_b \approx \frac{d_{min}}{n} P_w$$

These bounds and approximations were used in Fig. 2 to evaluate the performance of the (60,36) binary code derived from the (15,9) RS code with $b = 0$.

At a high signal-to-noise ratio (SNR), the approximation $\operatorname{erfc}(x) \approx e^{-x^2}/x\sqrt{\pi}$ may be used. Considering only the contribution of codewords at d_{min} , for soft decoding, one has the approximation

$$P_w \approx \frac{1}{2} A_{d_{min}} \frac{e^{-u^2}}{u\sqrt{\pi}}$$

where $u = \sqrt{Rd_{min}E_b/N_o}$. The probability of bit error P_b may be approximated by $P_b \approx (d_{min}/n)P_w$, as shown in Fig. 2.

Experience with simulation results for smaller codes indicates that this approximation is usually close to the true performance, while the bounds become loose at P_b larger than 10^{-6} .

IV. Conclusion

By computing the binary weight distribution of block codes, it is possible to estimate their performance with soft, maximum-likelihood decoding. This is useful in order to find large binary block codes with good performance, and to verify the performance obtained by specific soft-decoding algorithms presently under development.

References

- [1] R. Blahut, *Theory and Practice of Error Control Codes*, Reading, Massachusetts: Addison-Wesley, 1983.
- [2] R. J. McEliece, *The Theory of Information and Coding*, Reading, Massachusetts: Addison-Wesley, 1977.
- [3] F. J. MacWilliams and N. J. Sloane, *The Theory of Error-Correcting Codes*, New York: North-Holland Publishing Co., 1977.
- [4] T. Verhoeff, "An Updated Table of Minimum-Distance Bounds for Binary Linear Codes," *IEEE Transactions on Information Theory*, vol. IT-33, no. 5, pp. 65-80, September 1987.
- [5] A. J. Viterbi and J. K. Omura, *Principles of Digital Communication and Coding*, New York: McGraw-Hill, 1979.

Table 1. Binary weight distributions for the (7,2) and (7,5) codes.

weight	(21,6) CODE			(21,15) CODE		
	b=0, b=1	b=2, b=6	b=3, b=4, b=5	b=0, b=4	b=1, b=2, b=3	b=5, b=6
0	1	1	1	1	1	1
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	0	0	0	28	21	0
4	0	0	0	84	91	210
5	0	0	0	273	322	0
6	0	0	0	924	875	1638
7	3	0	0	1956	1809	0
8	0	21	14	2982	3129	6468
9	7	0	0	4340	4585	0
10	21	0	21	5796	5551	10878
11	21	0	0	5796	5551	0
12	7	42	21	4340	4585	9310
13	0	0	0	2982	3129	0
14	3	0	7	1956	1809	3570
15	0	0	0	924	875	0
16	0	0	0	273	322	651
17	0	0	0	84	91	0
18	0	0	0	28	21	42
19	0	0	0	0	0	0
20	0	0	0	0	0	0
21	1	0	0	1	1	0

Table 2. Weight distributions of the (60,24) code.

weight	b=0, b=5	b=1, b=4	b=2, b=3	b=6, b=14	b=7, b=13	b=8, b=12	b=9, b=11	b = 10
0	1	1	1	1	1	1	1	1
1	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	12
11	0	0	0	0	0	0	0	0
12	0	0	15	30	30	0	15	0
13	15	75	90	0	0	0	0	0
14	150	300	180	450	375	420	390	465
15	676	859	679	0	0	0	0	0
16	2250	2160	2490	5190	4125	4530	4500	4425
17	6555	5520	5505	0	0	0	0	0
18	14720	13220	13265	23420	28760	27485	27225	27240
19	29565	29760	29955	0	0	0	0	0
20	56304	60690	60795	135420	120585	121875	123000	120204
21	113255	115460	117455	0	0	0	0	0
22	218760	206520	205410	361140	408810	407565	407565	416895
23	342285	342180	339525	0	0	0	0	0
24	493400	531470	525185	1185680	1058015	1060295	1056500	1043975
25	758583	756000	753105	0	0	0	0	0
26	1079040	1000860	1018335	1778220	2016660	2016945	2020005	2034210
27	1277425	1275280	1281835	0	0	0	0	0
28	1414125	1519215	1509690	3387720	3046005	3040095	3043830	3017910
29	1665945	1669170	1666155	0	0	0	0	0
30	1831108	1719736	1717876	3013272	3414132	3418617	3413237	3450383
31	1665945	1669170	1666155	0	0	0	0	0
32	1414125	1519215	1509690	3403485	3040170	3041160	3041205	3012720
33	1277425	1275280	1281835	0	0	0	0	0
34	1079040	1000860	1018335	1779060	2015760	2015895	2018235	2027940
35	758583	756000	753105	0	0	0	0	0
36	493400	531470	525185	1176580	1061395	1059385	1058160	1057440
37	342285	342180	339525	0	0	0	0	0
38	218760	206520	205410	360300	409950	408615	409575	411105
39	113255	115460	117455	0	0	0	0	0
40	56304	60690	60795	138168	119493	122493	122148	119361
41	29565	29760	29955	0	0	0	0	0
42	14720	13220	13265	23780	28100	27035	26375	28070
43	6555	5520	5505	0	0	0	0	0

Table 2 (contd).

weight	b=0, b=5	b=1, b=4	b=2, b=3	b=6, b=14	b=7, b=13	b=8, b=12	b=9, b=11	b = 10
44	2250	2160	2490	4890	4305	4245	4755	4350
45	676	859	679	0	0	0	0	0
46	150	300	180	390	525	495	465	480
47	15	75	90	0	0	0	0	0
48	0	0	15	20	20	65	30	30
49	0	0	0	0	0	0	0	0
50	0	0	0	0	0	0	0	0
51	0	0	0	0	0	0	0	0
52	0	0	0	0	0	0	0	0
53	0	0	0	0	0	0	0	0
54	0	0	0	0	0	0	0	0
55	0	0	0	0	0	0	0	0
56	0	0	0	0	0	0	0	0
57	0	0	0	0	0	0	0	0
58	0	0	0	0	0	0	0	0
59	0	0	0	0	0	0	0	0
60	1	1	1	0	0	0	0	0

Table 3. Weight distributions of the (60,36) code.

weight	b=0, b=5	b=1, b=4	b=2, b=3	b=6, b=14	b=7, b=8, b=12	b=9, b=11	b = 10	b = 13
0	1	1	1	1	1	1	1	1
1	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0
7	0	0	0	60	0	0	15	0
8	105	360	270	105	105	75	135	60
9	0	0	0	660	765	945	1065	1005
10	9135	8067	9012	4350	4470	4500	4380	4605
11	0	0	0	20940	21045	20655	19995	20505
12	171290	170045	166730	84250	84370	84360	85950	84955
13	0	0	0	307620	308790	306720	310305	306690
14	2051130	2063850	2069655	1036980	1029780	1033080	1025820	1025910
15	0	0	0	3169396	3166006	3172656	3163509	3171106
16	17857290	17841435	17827110	8879100	8926260	8909250	8920440	8933025
17	0	0	0	23084220	23077425	23080395	23067975	23087925
18	110247955	110242255	110291800	55357350	55138110	55169540	55153100	55148985
19	0	0	0	121876260	121900185	121870485	121962285	121868505
20	499868640	499744149	499677249	248880309	249779349	249773439	249831315	249692244

Table 3 (contd).

weight	b=0, b=5	b=1, b=4	b=2, b=3	b=6, b=14	b=7, b=8, b=12	b=9, b=11	b = 10	b = 13
21	0	0	0	475905260	475911560	475934000	475793915	475896440
22	1686545400	1687429560	1687309875	846944880	843913200	843841440	843707280	844133640
23	0	0	0	1393888920	1393820040	1393856400	1393933350	1393917240
24	4299960090	4297337520	4297910160	2140496050	2148328210	2148448550	2148756230	2148052240
25	0	0	0	3094399368	3094425258	3094374858	3094295130	3094397658
26	8326857870	8331803670	8330907000	4181824860	4166545740	4166496360	4165579800	4166607690
27	0	0	0	5245474360	5245577050	5245564790	5245776110	5245426450
28	12370476540	12363639450	12364329360	6158040345	6180719145	6180621765	6182602665	6181074915
29	0	0	0	6821742120	6821661060	6821687280	6821545530	6821775660
30	14091448412	14098870268	14098595918	7076641208	7050800888	7050973648	7048404136	7050221828
31	0	0	0	6821742120	6821661060	6821687280	6821545530	6821775660
32	12370288365	12363796815	12364040385	6158040345	6180719145	6180621765	6182602665	6181074915
33	0	0	0	5245474360	5245577050	5245564790	5245776110	5245426450
34	8327053230	8331595110	8331101280	4181824860	4166545740	4166496360	4165579800	4166607690
35	0	0	0	3094399368	3094425258	3094374858	3094295130	3094397658
36	4299922280	4297446770	4297957910	2140496050	2148328210	2148448550	2148756230	2148052240
37	0	0	0	1393888920	1393820040	1393856400	1393933350	1393917240
38	1686443640	1687462200	1687212030	846944880	843913200	843841440	843707280	844133640
39	0	0	0	475905260	475911560	475934000	475793915	475896440
40	499970973	499664856	499699626	248880309	249779349	249773439	249831315	249692244
41	0	0	0	121876260	121900185	121870485	121962285	121868505
42	110224195	110285095	110300020	55357350	55138110	55169540	55153100	55148985
43	0	0	0	23084220	23077425	23080395	23067975	23087925
44	17833530	17831625	17829870	8879100	8926260	8909250	8920440	8933025
45	0	0	0	3169396	3166006	3172656	3163509	3171106
46	2071290	2066730	2063835	1036980	1029780	1033080	1025820	1025910
47	0	0	0	307620	308790	306720	310305	306690
48	166120	167845	168400	84250	84370	84360	85950	84955
49	0	0	0	20940	21045	20655	19995	20505
50	8895	8715	9000	4350	4470	4500	4380	4605
51	0	0	0	660	765	945	1065	1005
52	360	345	225	105	105	75	135	60
53	0	0	0	60	0	0	15	0
54	0	0	15	0	0	0	0	0
55	0	0	0	0	0	0	0	0
56	0	0	0	0	0	0	0	0
57	0	0	0	0	0	0	0	0
58	0	0	0	0	0	0	0	0
59	0	0	0	0	0	0	0	0
60	0	0	0	1	1	1	1	1

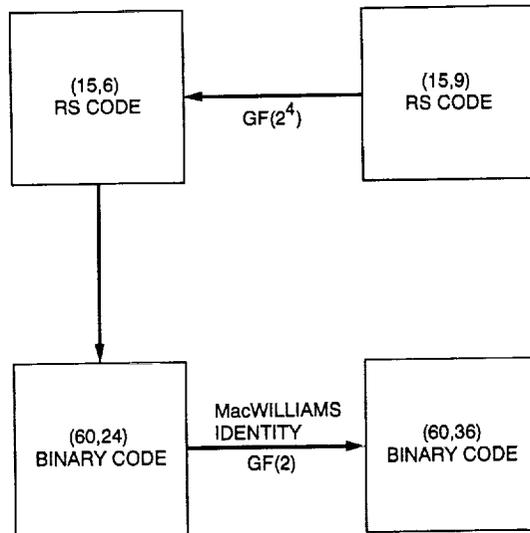


Fig. 1. Method used to find the binary weight distribution.

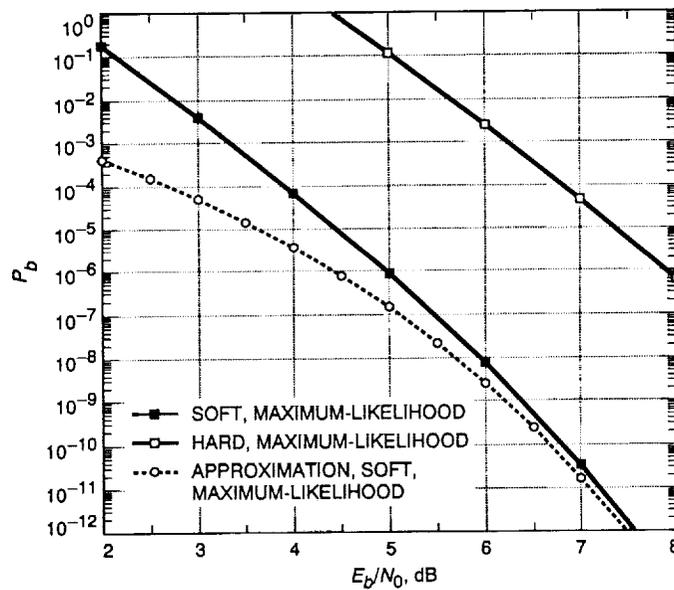


Fig. 2. Performance of (60,36) binary code derived from (15,9) RS code.

518 - ~~27~~ - 19431
 128451
 p-15

Multiple Symbol Partially Coherent Detection of MPSK

M. K. Simon

Telecommunications Systems Section

D. Divsalar

Communications Systems Research Section

In this article, it is shown that by using the known (or estimated) value of carrier tracking loop SNR in the decision metric, it is possible to improve the error probability performance of a partially coherent multiple phase-shift-keying (MPSK) system relative to that corresponding to the commonly used ideal coherent decision rule. Using a maximum-likelihood approach, an optimum decision metric is derived and shown to take the form of a weighted sum of the ideal coherent decision metric (i.e., correlation) and the noncoherent decision metric which is optimum for differential detection of MPSK. The performance of a receiver based on this optimum decision rule is derived and shown to provide continued improvement with increasing length of observation interval (data symbol sequence length). Unfortunately, increasing the observation length does not eliminate the error floor associated with the finite loop SNR. Nevertheless, in the limit of infinite observation length, the average error probability performance approaches the algebraic sum of the error floor and the performance of ideal coherent detection, i.e., at any error probability above the error floor, there is no degradation due to the partial coherence. It is shown that this limiting behavior is virtually achievable with practical size observation lengths. Furthermore, the performance is quite insensitive to mismatch between the estimate of loop SNR (e.g., obtained from measurement) fed to the decision metric and its true value. These results may be of use in low-cost Earth-orbiting or deep-space missions employing coded modulations.

I. Introduction

It is well known that for ideal phase coherent detection of multiple phase-shift-keying (MPSK), the decision rule that minimizes average bit error probability is based on a correlation metric and leads to bit-by-bit decisions. In practical situations, the phase introduced by the transmission over the channel is unknown and thus the assumption of perfect knowledge of this parameter at the receiver is idealistic. Typically, if the channel phase is reasonably well behaved, the receiver will attempt to estimate it via

some type of phase synchronization subsystem, such as a carrier phase tracking loop. Since the estimate is made in the presence of the ever-present additive channel thermal noise, the receiver's phase estimate used for demodulation purposes is not perfect. Detection under these circumstances is known as *partially coherent* detection.

Ordinarily in this environment, one continues to use the ideal coherent detection correlation metric despite the fact that it is no longer optimum for partially coherent detec-

tion. In particular, the presence of a phase error between the true channel and the receiver's estimate of it introduces memory into the observation, and thus any metric leading to bit-by-bit detection cannot be optimum. Instead, one must resort to sequence estimation where the length of the sequence is proportional to the duration over which the phase error can be assumed constant.

In this article, a maximum-likelihood approach to partially coherent detection is taken, an approach not unlike that previously applied to noncoherent and coherent detection. It will be shown that considerable performance improvement can be gained by using the optimum metric which leads to a maximum-likelihood sequence estimation (MLSE) type of algorithm.

II. Maximum-Likelihood Partially Coherent Detection of MPSK Over an AWGN Channel

Consider the transmission of MPSK signals over an additive white Gaussian noise (AWGN) channel. The baseband representation of the transmitted signal in the interval $(kT, (k+1)T)$ has the complex form

$$s_k = \sqrt{2P}e^{j\phi_k} \quad (1)$$

where P denotes the constant signal power, T denotes the MPSK symbol interval, and ϕ_k the transmitted phase which takes on one of M uniformly distributed values $\beta_m = 2\pi m/M; m = 0, 1, \dots, M-1$ around the unit circle. Assume that in addition to AWGN, the channel introduces a phase θ which can be constant (independent of time) over a duration of N data symbols and uniformly distributed in the interval $(-\pi, \pi)$. Thus, the received sequence \mathbf{r} is expressed as

$$\mathbf{r} = \mathbf{s}e^{j\theta} + \mathbf{n} \quad (2)$$

where $\mathbf{r} = (r_0, r_1, \dots, r_{N-1})$, $\mathbf{s} = (s_0, s_1, \dots, s_{N-1})$, and $\mathbf{n} = (n_0, n_1, \dots, n_{N-1})$ are the received sequence, transmitted sequence, and noise sequence, respectively. Also, n_k is a sample of zero mean complex Gaussian noise with variance (per dimension) $\sigma_n^2 = N_0/T$ where N_0 is the single-sided power spectral density of the noise process $n(t)$ at the receiver input.

For partially coherent detection, the receiver provides a carrier phase synchronization subsystem, e.g., a tracking loop, which derives a complex reference signal $e^{j\hat{\theta}}$ whose phase $\hat{\theta}$ is an estimate of the unknown channel phase θ . After demodulating \mathbf{r} with this reference (complex conjugate multiplication of the two signals), one gets

$$\mathbf{R} = \mathbf{r}e^{-j\hat{\theta}} = \mathbf{s}e^{j\phi_c} + \mathbf{n}e^{-j\hat{\theta}} \quad (3)$$

where $\phi_c \triangleq \theta - \hat{\theta}$ is the carrier phase error and typically has a Tikhonov probability density function (pdf) [1], i.e.,

$$p(\phi_c) = \frac{\exp(\rho \cos \phi_c)}{2\pi I_0(\rho)}; \quad |\phi_c| \leq \pi \quad (4)$$

Here ρ is a parameter related¹ to the tracking loop SNR and $I_0(\cdot)$ is the zeroth-order modified Bessel function of the first kind.

For the assumed AWGN model for \mathbf{n} , the a posteriori probability of the demodulated received sequence \mathbf{R} given the transmitted sequence \mathbf{s} and the carrier phase error ϕ_c follows from Eq. (3) and is

¹ For first-order tracking loops, ρ is indeed the loop SNR. For second-order loops, ρ is approximately the loop SNR for sufficiently large values [1]. In what follows, ρ is referred to simply as the loop SNR which is assumed to be known or estimated.

$$\begin{aligned} p(\mathbf{R}|\mathbf{s}, \phi_c) &= \frac{1}{(2\pi\sigma_n^2)^N} \exp \left\{ -\frac{\|\mathbf{R} - \mathbf{s}e^{j\phi_c}\|^2}{2\sigma_n^2} \right\} \\ &= \frac{1}{(2\pi\sigma_n^2)^N} \exp \left\{ -\frac{1}{2\sigma_n^2} \left[\sum_{i=0}^{N-1} [|R_{k-i}|^2 + |s_{k-i}|^2] - 2 \left| \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right| \cos(\phi_c - \alpha) \right] \right\} \end{aligned} \quad (5)$$

where

$$\alpha = \tan^{-1} \frac{\operatorname{Im} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\}}{\operatorname{Re} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\}} \quad (6)$$

Averaging Eq. (5) over the pdf in Eq. (4) gives, upon simplification,

$$\begin{aligned} p(\mathbf{R}|\mathbf{s}) &= \int_{-\pi}^{\pi} p(\mathbf{R}|\mathbf{s}, \phi_c) p(\phi_c) d\phi_c \\ &= \frac{1}{I_0(\rho)} \frac{1}{(2\pi\sigma_n^2)^N} \exp \left\{ -\frac{1}{2\sigma_n^2} \sum_{i=0}^{N-1} [|R_{k-i}|^2 + |s_{k-i}|^2] \right\} \\ &\quad \times I_0 \left(\frac{1}{\sigma_n^2} \sqrt{ \left(\left| \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right| \cos \alpha + \rho\sigma_n^2 \right)^2 + \left(\left| \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right| \sin \alpha \right)^2 } \right) \end{aligned} \quad (7)$$

Since

$$\left| \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right| \cos \alpha = \operatorname{Re} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\}; \quad \left| \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right| \sin \alpha = \operatorname{Im} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\} \quad (8)$$

Eq. (7) further simplifies to

$$\begin{aligned} p(\mathbf{R}|\mathbf{s}) &= \frac{1}{I_0(\rho)} \frac{1}{(2\pi\sigma_n^2)^N} \exp \left\{ -\frac{1}{2\sigma_n^2} \sum_{i=0}^{N-1} [|R_{k-i}|^2 + |s_{k-i}|^2] \right\} \\ &\quad \times I_0 \left(\frac{1}{\sigma_n^2} \sqrt{ \left(\operatorname{Re} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\} + \rho\sigma_n^2 \right)^2 + \left(\operatorname{Im} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\} \right)^2 } \right) \end{aligned} \quad (9)$$

Note from Eq. (1) that for MPSK, $|s_k|^2$ is constant for all transmitted phases β_m . Thus, since $I_0(x)$ is a monotonic function of its argument, maximizing $p(\mathbf{R}|\mathbf{s})$ over \mathbf{s} is equivalent to finding

$$\begin{aligned} \max_{\mathbf{s}} \left\{ \left(\operatorname{Re} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\} + \rho\sigma_n^2 \right)^2 + \left(\operatorname{Im} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\} \right)^2 \right\} = \\ \max_{\mathbf{s}} \left\{ \left| \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* + \rho\sigma_n^2 \right|^2 \right\} = \max_{\mathbf{s}} \left\{ \left| \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right|^2 + 2\rho\sigma_n^2 \operatorname{Re} \left\{ \sum_{i=0}^{N-1} R_{k-i} s_{k-i}^* \right\} \right\} \end{aligned} \quad (10)$$

This, using Eq. (1), results in the decision rule

$$\text{choose } \hat{\phi}_k, \hat{\phi}_{k-1}, \dots, \hat{\phi}_{k-N+1} \text{ if } \left\{ \left| \sum_{i=0}^{N-1} R_{k-i} e^{-j\hat{\phi}_{k-i}} \right|^2 + \frac{2\rho\sigma_n^2}{\sqrt{2P}} \operatorname{Re} \left\{ \sum_{i=0}^{N-1} R_{k-i} e^{-j\hat{\phi}_{k-i}} \right\} \right\} \text{ is maximum} \quad (11)$$

where $\hat{\phi}_k, \hat{\phi}_{k-1}, \dots, \hat{\phi}_{k-N+1}$ is a particular sequence of the transmitted phases β_m . In Eq. (11), the first term inside the braces represents the component of the decision metric associated with *noncoherent* (differential) detection [2], i.e., total lack of knowledge of the uniformly distributed channel phase θ . The second term inside the braces represents the component of the decision metric for ideal *coherent* detection, i.e., complete knowledge of the channel phase θ . Thus, the partially coherent decision metric is a linear combination of the coherent and noncoherent decision metrics with the weighting of the two terms in proportion to the product of the tracking loop SNR and the channel noise variance. Note that for any nonzero value of ρ , this decision rule is unique because the second term inside the braces in Eq. (11) is unique but not the first term. For $\rho = 0$, which corresponds to differentially coherent detection, there is a phase ambiguity since the addition of an arbitrary fixed phase, say ϕ_a , to all N estimated phases $\hat{\phi}_k, \hat{\phi}_{k-1}, \dots, \hat{\phi}_{k-N+1}$ results in the same decision for ϕ . In [2], the authors observed that by letting $\phi_a = \phi_{k-N+1}$ and differentially encoding the input phases at the transmitter,

$$\phi_k = \phi_{k-1} + \Delta\phi_k \quad (12)$$

where now $\Delta\phi_k$ denotes the input data phase corresponding to the k th transmission interval and ϕ_k the differentially encoded version of it, the decision rule can turn into one in which the phase ambiguity is resolved. From now on, assume $\rho \neq 0$ and thus that there is no formal requirement for differentially encoding the data phase symbols.

Figure 1 is an illustration in complex form of a receiver implemented on the basis of Eq. (11). Note that this receiver requires knowledge of the loop SNR ρ , the signal power P , and the noise variance σ_n^2 . The accuracy of this knowledge, which must be obtained by measurement, will have an impact on the ultimate performance of this receiver. Later, in Subsection E, the authors investigate the sensitivity of the receiver to a mismatch between the true loop SNR and the value supplied to the receiver implementation in Fig. 1. In the next section, except in Section III.E, it is assumed that the receiver

has perfect knowledge of ρ , and thus should outperform a conventional bit-by-bit correlation receiver which does not make use of this knowledge. The following sections determine how much the optimum partially coherent sequence receiver outperforms the conventional bit-by-bit correlation receiver.

III. Bit Error Probability Performance

To obtain a simple upper bound on the average bit error probability, P_b , of the proposed N -bit detection scheme, use a union bound analogous to that used for upper bounding the performance of error correction coded systems. In particular, the upper bound on P_b is the sum of the pairwise error probabilities associated with each N -bit error sequence. Each pairwise error probability is then either evaluated directly or itself upper bounded. Mathematically speaking, let $\phi = (\phi_k, \phi_{k-1}, \dots, \phi_{k-N+1})$ denote the sequence of N transmitted information phases and $\hat{\phi} = (\hat{\phi}_k, \hat{\phi}_{k-1}, \dots, \hat{\phi}_{k-N+1})$ be the corresponding sequence of detected phases. Let \mathbf{u} be the sequence of $b = N \log_2 M$ information bits that produces ϕ at the transmitter and let $\hat{\mathbf{u}}$ be the sequence of b bits that results from the detection of $\hat{\phi}$. Then, since MPSK is a symmetric signalling set, i.e., it satisfies a uniform error probability (UEP) criterion, one gets an upper bound on the bit error probability,

$$P_b(\phi_c) \leq \frac{1}{N \log_2 M} \sum_{\hat{\phi} \neq \phi} w(\mathbf{u}, \hat{\mathbf{u}}) \operatorname{Pr} \{ \hat{\eta} > \eta | \phi, \phi_c \} \quad (13)$$

where the decision statistic η is defined from Eqs. (10) and (11) by²

$$\eta = \left| \sum_{i=0}^{N-1} R_{k-i} e^{-j\phi_{k-i}} + \frac{\rho\sigma_n^2}{\sqrt{2P}} \right|^2 \quad (14)$$

² Note that when compared with Eq. (11), η of Eq. (14) includes the additional constant $(\rho\sigma_n^2/\sqrt{2P})^2$. This, however, has no effect on the decision-making process and thus one can use the convenient form of Eq. (14) with no loss in generality.

and the corresponding error statistic $\hat{\eta}$ is identical to Eq. (14) with each ϕ_k replaced by $\hat{\phi}_k$. In Eq. (13), $w(\mathbf{u}, \hat{\mathbf{u}})$ denotes the Hamming distance between \mathbf{u} and $\hat{\mathbf{u}}$, ϕ is any input sequence (e.g., the null sequence $(0, 0, \dots, 0) = 0$), and $\Pr\{\hat{\eta} > \eta \mid \phi, \phi_c\}$ denotes the pairwise probability that $\hat{\phi}$ is incorrectly chosen when indeed ϕ was sent. Note that the bound in Eq. (13) is computed for a fixed carrier phase error, ϕ_c , which accounts for the notational dependence of $\Pr\{\hat{\eta} > \eta \mid \phi, \phi_c\}$ and thus $P_b(\phi_c)$ on ϕ_c .

ity that $\hat{\phi}$ is incorrectly chosen when indeed ϕ was sent. Note that the bound in Eq. (13) is computed for a fixed carrier phase error, ϕ_c , which accounts for the notational dependence of $\Pr\{\hat{\eta} > \eta \mid \phi, \phi_c\}$ and thus $P_b(\phi_c)$ on ϕ_c .

A. Evaluation of the Pairwise Error Probability

To compute $\Pr\{\hat{\eta} > \eta \mid \phi, \phi_c\}$, the approach taken in [2] is used for evaluating the performance of multiple symbol differentially coherent detection of MPSK. In particular, letting $\eta = |z_1|^2$ and $\hat{\eta} = |z_2|^2$ [see Eq. (14) and the statement below it for the definitions of z_1 and z_2], then [3]

$$\Pr\{\hat{\eta} > \eta \mid \phi, \phi_c\} = \frac{1}{2} \left[1 - Q(\sqrt{b}, \sqrt{a}) + Q(\sqrt{a}, \sqrt{b}) \right] \triangleq f(a, b) \quad (15)$$

where $Q(x, y)$ is the Marcum Q function [4] and

$$\left\{ \begin{array}{l} b \\ a \end{array} \right\} = \frac{1}{2N_z} \left\{ \frac{S_1 + S_2 - 2|\xi| \sqrt{S_1 S_2} \cos(\theta_1 - \theta_2 + \nu)}{1 - |\xi|^2} \pm \frac{S_1 - S_2}{\sqrt{1 - |\xi|^2}} \right\} \quad (16)$$

where the + sign and - sign correspond to b and a , respectively, and

$$S_1 = P \left| N + \frac{\rho}{2E_s/N_0} e^{-j\phi_c} \right|^2 = P \left(N^2 + \frac{\rho}{E_s/N_0} N \cos \phi_c + \left(\frac{\rho}{2E_s/N_0} \right)^2 \right)$$

$$S_2 = P \left| \delta + \frac{\rho}{2E_s/N_0} e^{-j\phi_c} \right|^2 = P \left(|\delta|^2 + \frac{\rho}{E_s/N_0} \operatorname{Re} \{ \delta e^{j\phi_c} \} + \left(\frac{\rho}{2E_s/N_0} \right)^2 \right)$$

$$N_z = \frac{1}{2} |z_1 - \bar{z}_1|^2 = N \frac{N_0}{T}$$

$$\xi = \frac{1}{2N_z} \overline{(z_1 - \bar{z}_1)} (z_2 - \bar{z}_2)^* = \frac{\delta}{N}; \quad \nu = \arg \xi = \arg \delta$$

$$\theta_1 = \arg \bar{z}_1 = \arg \left\{ N e^{j\phi_c} + \frac{\rho}{2E_s/N_0} \right\}; \quad \theta_2 = \arg \bar{z}_2 = \arg \left\{ \delta e^{j\phi_c} + \frac{\rho}{2E_s/N_0} \right\} \quad (17)$$

and

$$\delta = \sum_{i=0}^{N-1} e^{j(\phi_{k-i} - \hat{\phi}_{k-i})} \quad (18)$$

which is a normalized time cross correlation between \mathbf{s} and $\hat{\mathbf{s}}$. Also, $E_s/N_0 \triangleq PT/N_0$ denotes the symbol energy-to-noise spectral density ratio and is related to the bit energy-to-noise spectral density ratio E_b/N_0 by $E_s/N_0 = (E_b/N_0) \log_2 M$. Substituting Eq. (17) into Eq. (16) results, after considerable simplification, in

$$\left\{ \begin{array}{l} b \\ a \end{array} \right\} = \frac{E_s}{2N_0} \left\{ N \left[1 + \frac{1}{N} \left(\frac{\rho}{E_s/N_0} \right) \cos \phi_c + \frac{1}{2N(N^2 - |\delta|^2)} \left(\frac{\rho}{E_s/N_0} \right)^2 (N - |\delta| \cos \nu) \right] \right. \\ \left. \pm \frac{E_s}{2N_0} \left[\sqrt{N^2 - |\delta|^2} + \frac{1}{\sqrt{N^2 - |\delta|^2}} \left(\frac{\rho}{E_s/N_0} \right) (N \cos \phi_c - |\delta| \cos(\phi_c + \nu)) \right] \right\} \quad (19)$$

Now some special cases of practical interest are considered.

B. Case 1: Binary PSK With Two-Symbol Observation and Detection ($M = 2, N = 2$)

In this case, $E_s/N_0 = E_b/N_0$. There are $M^2 - 1 = 3$ possible error sequences each of length 2. The pertinent results related to the evaluation of Eqs. (18) and (19) are

$\phi_k - \hat{\phi}_k$	$\phi_{k-1} - \hat{\phi}_{k-1}$	δ
0	π	0
π	0	0
π	π	-2

For the first two error sequences, Eq. (19) evaluates to

$$b = \frac{E_b}{2N_0} \left[4 + 2 \left(\frac{\rho}{E_b/N_0} \right) \cos \phi_c + \frac{1}{4} \left(\frac{\rho}{E_b/N_0} \right)^2 \right] \\ a = \frac{E_b}{2N_0} \left[\frac{1}{4} \left(\frac{\rho}{E_b/N_0} \right)^2 \right] \quad (20)$$

For the third error sequence, both a and b approach infinity (the ratio a/b , however, approaches unity) as δ approaches -2 . Thus, one must evaluate the pairwise error probability Eq. (15) separately for this case. It is straightforward to show that

$$\lim_{\substack{a \rightarrow \infty \\ b \rightarrow \infty \\ a/b \rightarrow 1}} f(a, b) = \frac{1}{2} \operatorname{erfc} \left(\sqrt{\frac{b}{2}} - \sqrt{\frac{a}{2}} \right) \quad (21)$$

Furthermore, in the general case where $\delta \rightarrow -N$, Eq. (21) evaluates to

$$\lim_{\delta \rightarrow -N} f(a, b) = \frac{1}{2} \operatorname{erfc} \left(\sqrt{N \frac{E_b}{N_0}} \cos \phi_c \right) \quad (22)$$

which for $N = 2$ and $M = 2$ becomes

$$\lim_{\delta \rightarrow -2} f(a, b) = \frac{1}{2} \operatorname{erfc} \left(\sqrt{\frac{2E_b}{N_0}} \cos \phi_c \right) \quad (23)$$

Finally, noting that the Hamming distance $w(\mathbf{u}, \hat{\mathbf{u}})$ is equal to 1 for the first two error sequences and is equal to 2 for the third sequence, substituting Eqs. (23) and (15) combined with Eqs. (19) and (20) into the expression for bit error probability in Eq. (13) gives

$$P_b(\phi_c) \leq f \left(\frac{E_b}{N_0} \left[\frac{1}{8} \left(\frac{\rho}{E_b/N_0} \right)^2 \right], \frac{E_b}{N_0} \left[2 + \left(\frac{\rho}{E_b/N_0} \right) \cos \phi_c + \frac{1}{8} \left(\frac{\rho}{E_b/N_0} \right)^2 \right] \right) + \frac{1}{2} \operatorname{erfc} \left(\sqrt{\frac{2E_b}{N_0}} \cos \phi_c \right) \quad (24)$$

Finally, the upper bound on average bit error probability P_b is obtained by averaging the upper bound in Eq. (24) over the pdf in Eq. (4). Figures 2 and 3 are plots of this upper bound on average probability versus E_b/N_0 in decibels for values of $\rho = 7$ dB and 10 dB, respectively. For the purpose of comparison, the exact results (i.e., not an upper bound) for the conventional ideal coherent metric operating in a noisy carrier synchronization environment are

$$P_b = \int_{-\pi}^{\pi} \frac{1}{2} \operatorname{erfc} \left(\sqrt{\frac{E_b}{N_0}} \cos \phi_c \right) p(\phi_c) d\phi_c \quad (25)$$

where $p(\phi_c)$ is given by Eq. (4). Even with only one additional observation symbol interval, considerable savings in E_b/N_0 can be achieved at a fixed error probability, particularly in the region of the knee of the curve where the system begins approaching its irreducible error probability³ asymptote (error floor).

C. Case 2: Binary PSK with N -Symbol Observation and Detection ($M = 2$, N arbitrary)

For N arbitrary, δ takes on values $-(N - 2i); i = 0, 1, 2, \dots, N - 1$. The number of error sequences corresponding to each of these values of δ is binomially distributed, i.e., there are $\binom{N}{i}$ sequences that yield a value $\delta = -(N - 2i)$. Furthermore, the Hamming weight associated with each of the $\binom{N}{i}$ sequences that yield a value $\delta = -(N - 2i)$ is $w(\mathbf{u}, \hat{\mathbf{u}}) = N - i$. Finally then, using the above in Eqs. (19) and (22) and substituting the results in Eq. (13), the conditional bit error probability is upper bounded by

$$\begin{aligned} P_b(\phi_c) &\leq \frac{1}{N} \left[N \frac{1}{2} \operatorname{erfc} \left(\sqrt{N \frac{E_b}{N_0}} \cos \phi_c \right) + \sum_{i=1}^{N-1} \binom{N}{i} (N - i) f(a_i, b_i) \right] \\ &= \frac{1}{2} \operatorname{erfc} \left(\sqrt{N \frac{E_b}{N_0}} \cos \phi_c \right) + \sum_{i=1}^{N-1} \binom{N-1}{i} f(a_i, b_i) \end{aligned} \quad (26)$$

where

$$\left\{ \begin{array}{l} b_i \\ a_i \end{array} \right\} = \frac{E_b}{2N_0} \left\{ \left[N + \left(\frac{\rho}{E_b/N_0} \right) \cos \phi_c + \frac{1}{4i} \left(\frac{\rho}{2E_b/N_0} \right)^2 \right] \pm \left[2\sqrt{i(N-i)} + \sqrt{\frac{N-i}{i}} \left(\frac{\rho}{E_b/N_0} \right) \cos \phi_c \right] \right\} \quad (27)$$

³ It is well known [1] that conventional PSK systems exhibit an irreducible error probability (i.e., a finite error probability in the limit as E_b/N_0 approaches infinity) when a noisy carrier synchronization reference with fixed power is used as a demodulation signal. This is observed by examining a curve of P_b versus E_b/N_0 with loop SNR, ρ , held fixed. The value of this irreducible error probability is given by [1] $P_b|_{\text{irr}} = \int_{\rho/2}^{\rho} p(\phi_c) d\phi_c$. Note that in practice, as the observation length increases, one should decrease the loop bandwidth of the phase-locked loop (PLL), which results in an increase in the loop SNR. Also, as the bit SNR increases, the loop SNR (for fixed modulation index) increases and thus the error floor decreases.

The upper bound on unconditional average bit error probability is now obtained by averaging Eq. (26) over the pdf in Eq. (4). The numerical results are illustrated in Figs. 2 and 3 for values of $N = 4, 6,$ and 8 . As N gets large, the curves appear to approach an asymptote. This asymptotic behavior is analytically evaluated as follows:

For large N , the first term in Eq. (26) when integrated over the pdf in Eq. (4) approaches the irreducible error probability $P_b|_{irr} = \int_{\pi/2}^{\pi} p(\phi_c) d\phi_c$. Also, the dominant term in the summation term of Eq. (26) corresponds to $i = N-1$, i.e., $\delta = N-2$. Thus, for large N , the second term of Eq. (26) approaches $f(a_{N-1}, b_{N-1})$ where

$$\left\{ \begin{matrix} b_{N-1} \\ a_{N-1} \end{matrix} \right\} \cong \frac{E_b}{2N_0} \left\{ N \pm 2\sqrt{(N-1)} \right\} \quad (28)$$

Since from Eq. (28), $\sqrt{b_{N-1}} \gg \sqrt{b_{N-1}} - \sqrt{a_{N-1}}$, then using the asymptotic form of Eq. (15) for a and b large (see Appendix A of [5]), namely,

$$f(a, b) \cong \frac{1}{2} \operatorname{erfc} \left(\frac{\sqrt{b} - \sqrt{a}}{2} \right) \quad (29)$$

The value $f(a_{N-1}, b_{N-1})$ is obtained as

$$f(a_{N-1}, b_{N-1}) \cong \frac{1}{2} \operatorname{erfc} \left\{ \sqrt{\frac{E_b}{N_0}} \right\} \quad (30)$$

independent of ϕ_c . Finally then, for large N , the asymptotic behavior of the average bit error probability is approximately upper bounded by

$$P_b \leq \frac{1}{2} \operatorname{erfc} \sqrt{\frac{E_b}{N_0}} + 2 \int_{\pi/2}^{\pi} p(\phi_c) d\phi_c \quad (31)$$

namely, the sum of the bit error probability for ideal coherent detection and the error floor. Equation (31) is in very close agreement with the curves for $N = 8$ in Figs. 2 and 3.

D. Case 3: Quaternary PSK With Two-Symbol Observation and Detection ($M = 4, N = 2$)

In this case, $E_s/N_0 = 2E_b/N_0$. There are now a total of $M^2 - 1 = 15$ possible error sequences each of length 2. Of these, only eight produce distinct combinations of $|\delta|$ and ν . These are tabulated below:

Error sequence	Case	$\phi_k - \hat{\phi}_k$	$\phi_{k-1} - \hat{\phi}_{k-1}$	$ \delta $	ν
1,2,3,4	1	$\pi, 0, 3\pi/2, \pi/2$	$0, \pi, \pi/2, 3\pi/2$	0	0
5,6	2	$0, \pi/2$	$\pi/2, 0$	$\sqrt{2}$	$\pi/4$
7,8	3	$0, 3\pi/2$	$3\pi/2, 0$	$\sqrt{2}$	$-\pi/4$
9,10	4	$\pi/2, \pi$	$\pi, \pi/2$	$\sqrt{2}$	$3\pi/4$
11,12	5	$\pi, 3\pi/2$	$3\pi/2, \pi$	$\sqrt{2}$	$-3\pi/4$
13	6	$\pi/2$	$\pi/2$	2	$\pi/2$
14	7	$3\pi/2$	$3\pi/2$	2	$-\pi/2$
15	8	π	π	2	π

The corresponding values of a and b for each of the first five cases which correspond to $|\delta| \neq 2$ are given as follows:

$$\begin{aligned} \left\{ \begin{matrix} b \\ a \end{matrix} \right\}_1 &= \frac{E_b}{N_0} \left\{ \left[2 + \left(\frac{\rho}{2E_b/N_0} \right) \cos \phi_c + \frac{1}{4} \left(\frac{\rho}{2E_b/N_0} \right)^2 \right] \pm \left[2 + \left(\frac{\rho}{2E_b/N_0} \right) \cos \phi_c \right] \right\} \\ \left\{ \begin{matrix} b \\ a \end{matrix} \right\}_2 &= \frac{E_b}{N_0} \left\{ \left[2 + \left(\frac{\rho}{2E_b/N_0} \right) \cos \phi_c + \frac{1}{4} \left(\frac{\rho}{2E_b/N_0} \right)^2 \right] \pm \left[\sqrt{2} + \frac{1}{\sqrt{2}} \left(\frac{\rho}{2E_b/N_0} \right) (\cos \phi_c + \sin \phi_c) \right] \right\} \\ \left\{ \begin{matrix} b \\ a \end{matrix} \right\}_3 &= \frac{E_b}{N_0} \left\{ \left[2 + \left(\frac{\rho}{2E_b/N_0} \right) \cos \phi_c + \frac{1}{4} \left(\frac{\rho}{2E_b/N_0} \right)^2 \right] \pm \left[\sqrt{2} + \frac{1}{\sqrt{2}} \left(\frac{\rho}{2E_b/N_0} \right) (\cos \phi_c - \sin \phi_c) \right] \right\} \\ \left\{ \begin{matrix} b \\ a \end{matrix} \right\}_4 &= \frac{E_b}{N_0} \left\{ \left[2 + \left(\frac{\rho}{2E_b/N_0} \right) \cos \phi_c + \frac{3}{4} \left(\frac{\rho}{2E_b/N_0} \right)^2 \right] \pm \left[\sqrt{2} + \frac{1}{\sqrt{2}} \left(\frac{\rho}{2E_b/N_0} \right) (3 \cos \phi_c + \sin \phi_c) \right] \right\} \\ \left\{ \begin{matrix} b \\ a \end{matrix} \right\}_5 &= \frac{E_b}{N_0} \left\{ \left[2 + \left(\frac{\rho}{2E_b/N_0} \right) \cos \phi_c + \frac{3}{4} \left(\frac{\rho}{2E_b/N_0} \right)^2 \right] \pm \left[\sqrt{2} + \frac{1}{\sqrt{2}} \left(\frac{\rho}{2E_b/N_0} \right) (3 \cos \phi_c - \sin \phi_c) \right] \right\} \quad (32) \end{aligned}$$

For cases 6 and 7, the following is analogous to Eq. (22):

$$\lim_{\delta \rightarrow \pm jN} f(a, b) = \frac{1}{2} \operatorname{erfc} \left(\sqrt{\frac{NE_b}{2N_0}} (\cos \phi_c \pm \sin \phi_c) \right) \quad (33)$$

which for $N = 2$ and $M = 4$ becomes

$$\lim_{\delta \rightarrow \pm j2} f(a, b) = \frac{1}{2} \operatorname{erfc} \left(\sqrt{\frac{2E_b}{N_0}} (\cos \phi_c \pm \sin \phi_c) \right) \quad (34)$$

Finally, for case 8, Eq. (22) is used to obtain

$$\lim_{\delta \rightarrow -2} f(a, b) = \frac{1}{2} \operatorname{erfc} \left(\sqrt{\frac{4E_b}{N_0}} \cos \phi_c \right) \quad (35)$$

Evaluating the Hamming distances for the 15 error sequences and substituting the above results into the expression for the bit error probability bound in Eq. (13) gives

$$\begin{aligned} P_b(\phi_c) &\leq \frac{1}{4} \{ 6f(a_1, b_1) + 2f(a_2, b_2) + 2f(a_3, b_3) + 4f(a_4, b_4) + 4f(a_5, b_5) \} \\ &+ \frac{1}{4} \left\{ \operatorname{erfc} \left(\sqrt{\frac{2E_b}{N_0}} (\cos \phi_c + \sin \phi_c) \right) + \operatorname{erfc} \left(\sqrt{\frac{2E_b}{N_0}} (\cos \phi_c - \sin \phi_c) \right) + \operatorname{erfc} \left(\sqrt{\frac{4E_b}{N_0}} \cos \phi_c \right) \right\} \quad (36) \end{aligned}$$

Figures 4 and 5 are comparable to Figs. 2 and 3 for the $M = 4$ (QPSK) case. The analytical exact result corresponding to the ideal coherent metric operating in a noisy carrier synchronization environment is now [1]

$$P_b = \int_{-\pi}^{\pi} \frac{1}{4} \operatorname{erfc} \left(\sqrt{\frac{E_b}{N_0}} (\cos \phi_c + \sin \phi_c) \right) p(\phi_c) d\phi_c + \int_{-\pi}^{\pi} \frac{1}{4} \operatorname{erfc} \left(\sqrt{\frac{E_b}{N_0}} (\cos \phi_c - \sin \phi_c) \right) p(\phi_c) d\phi_c \quad (37)$$

Analysis and plots for larger values of N are not included here, but they would show further improvement as was true for the binary case.

E. Performance Sensitivity to Mismatch In Loop SNR

Here the authors investigate the sensitivity of the average bit error probability (in terms of its upper bound) of the MLSE receiver to a mismatch between the true loop SNR, ρ , and the estimate of it, $\hat{\rho}$, supplied to the implementation of Fig. 1. In particular, the authors evaluate, for the special cases of Sections III.B and III.C, the upper bound

$$P_b \leq \int_{-\pi}^{\pi} P_{bu}(\phi_c; \hat{\rho}) p(\phi_c) d\phi_c \quad (38)$$

where $P_{bu}(\phi_c; \hat{\rho})$ is given by the upper bound in Eq. (24) or Eq. (26) with ρ replaced by $\hat{\rho} = \rho [1 + (\hat{\rho} - \rho)/\rho] \triangleq \rho(1 + \epsilon)$ and $p(\phi_c)$ is as given by Eq. (4). Figures 6 and 7 are illustrations of Eq. (38) for $M = 2$, $\rho = 10$ dB, and $N = 2$ and 8, respectively, with fractional mismatch ϵ as a parameter. One observes that even with mismatches as much as 50 percent ($\epsilon = \pm 0.5$), there is negligible effect on the error probability performance. Thus, the authors

conclude that the MLSE receiver is quite insensitive to mismatch in the loop SNR.

IV. Conclusions

By making use of the known (or estimated) value of loop SNR in the decision metric, it is possible to improve the error probability performance of a partially coherent MPSK system relative to that corresponding to the commonly used ideal coherent decision rule. Using a maximum-likelihood approach, an optimum decision metric was derived and shown to take the form of a weighted sum of the ideal coherent decision metric (i.e., correlation) and the noncoherent decision metric previously shown to be optimum for differential detection of MPSK. The performance of a receiver based on this optimum decision rule improves with the increasing length of the observation interval (data symbol sequence length). Furthermore, the performance is quite insensitive to mismatch between the estimate of loop SNR (e.g., obtained from measurement) fed to the decision metric relative and its true value.

References

- [1] W. C. Lindsey and M. K. Simon, *Telecommunication Systems Engineering*, New York: Prentice-Hall, Inc., 1973.
- [2] D. Divsalar and M. K. Simon, "Multiple-Symbol Differential Detection of MPSK," *IEEE Transactions on Communications*, vol. 38, no. 3, pp. 30-308, March 1990.
- [3] S. Stein, "Unified Analysis of Certain Coherent and Noncoherent Binary Communications Systems," *IEEE Transactions on Information Theory*, vol. IT-10, pp. 43-51, January 1964.
- [4] J. Marcum, *Tables of Q Functions*, RAND Corporation Report M-339, Santa Monica, California: Rand Corporation, January 1950.
- [5] M. Schwartz, W. R. Bennett, and S. Stein, *Communication Systems and Techniques*, New York: McGraw-Hill, Inc., 1966.

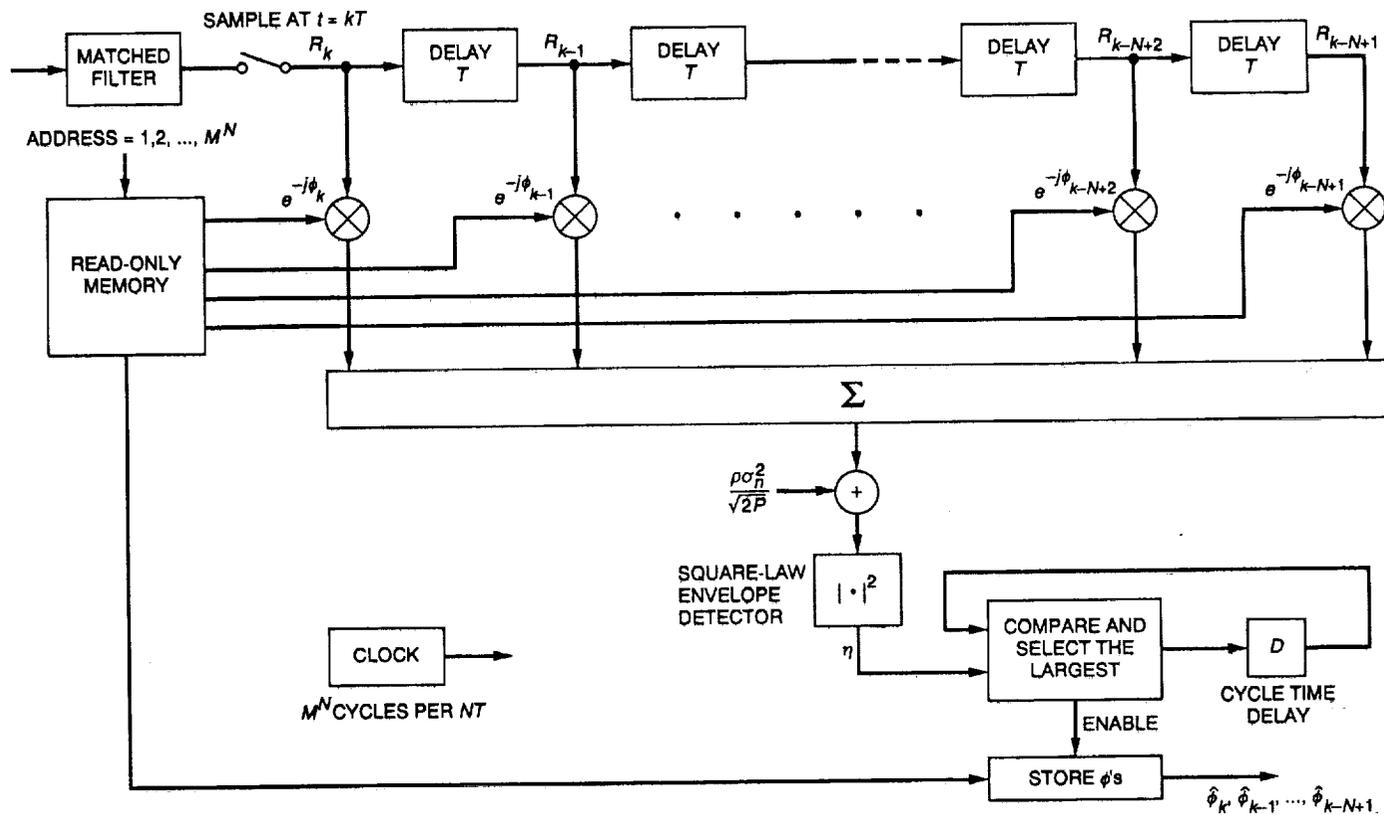


Fig. 1. Direct implementation of a receiver for multiple symbol detection of partially coherent MPSK.

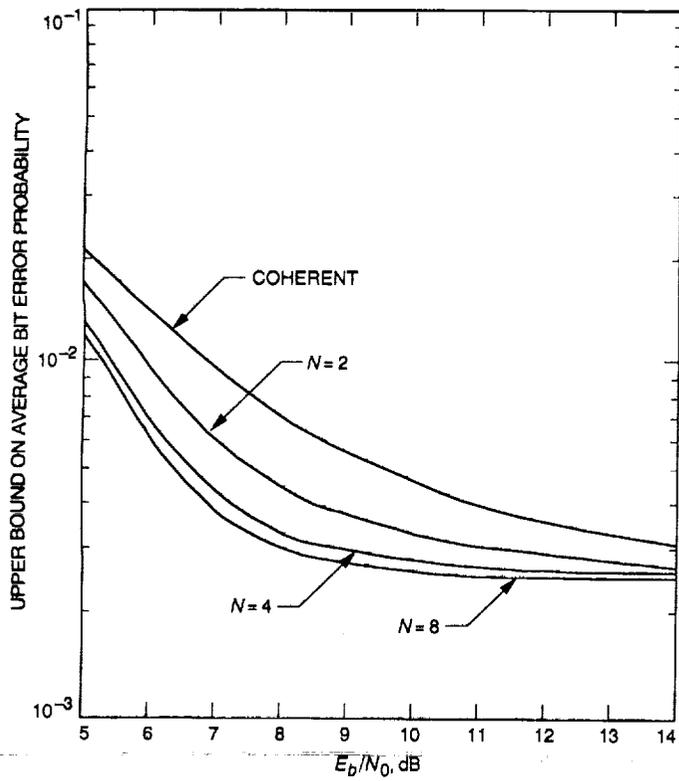


Fig. 2. Upper bound on average bit error probability versus E_b/N_0 in decibels for MLSE with N as a parameter; $M = 2$ and $\rho = 7$ dB.

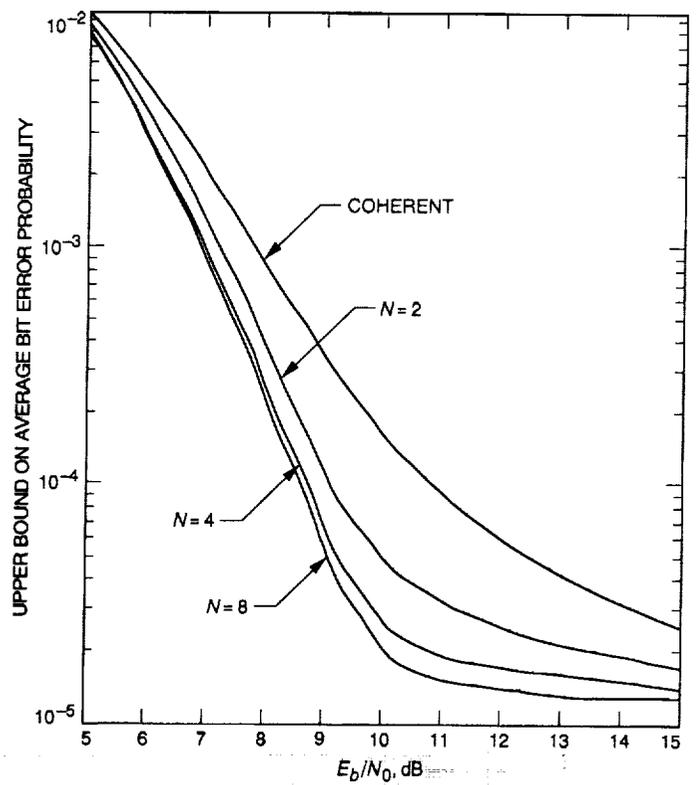


Fig. 3. Upper bound on average bit error probability versus E_b/N_0 in decibels for MLSE with N as a parameter; $M = 2$ and $\rho = 10$ dB.

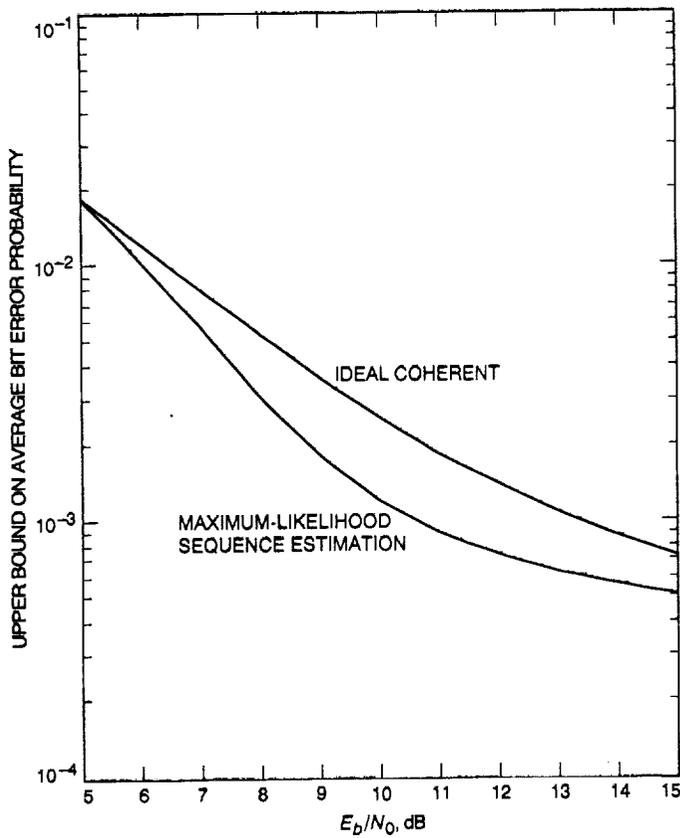


Fig. 4. Upper bound on average bit error probability versus E_b/N_0 in decibels for MLSE and comparison with exact performance of ideal coherent metric; $M = 4$, $N = 2$, and $\rho = 13$ dB.

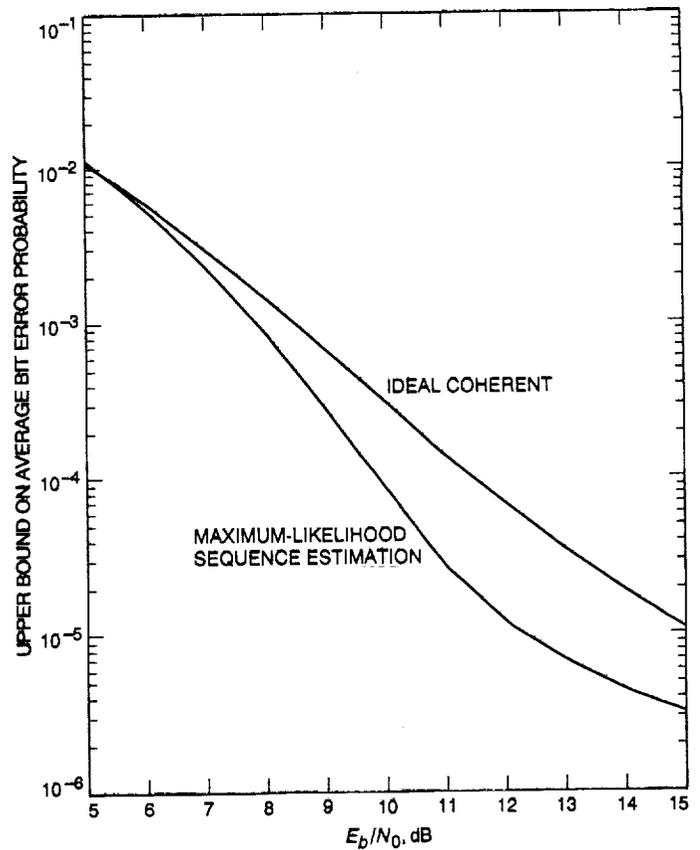


Fig. 5. Upper bound on average bit error probability versus E_b/N_0 in decibels for MLSE and comparison with exact performance of ideal coherent metric; $M = 4$, $N = 2$, and $\rho = 16$ dB.

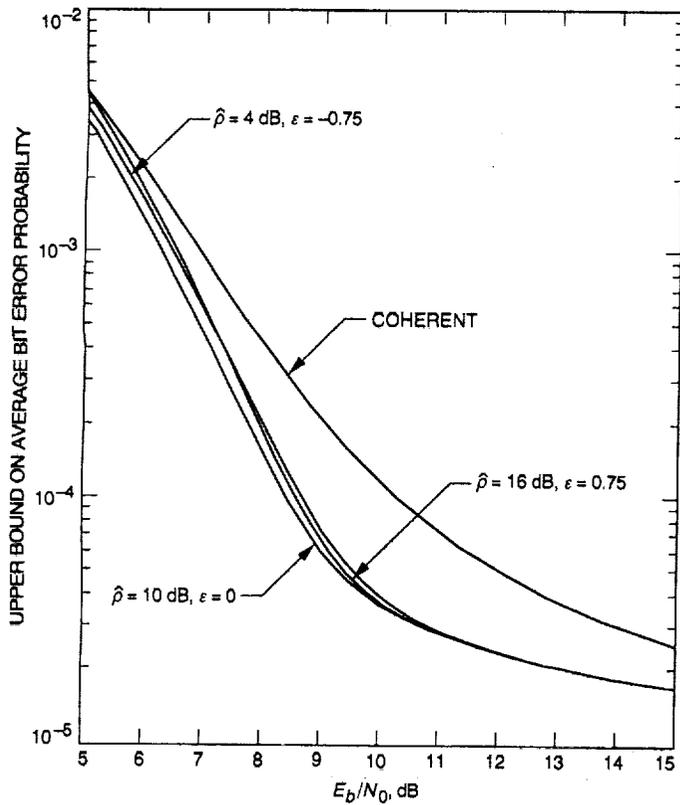


Fig. 6. Upper bound on average bit error probability versus E_b/N_0 in decibels for MLSE in the presence of loop SNR mismatch; $M = 2$, $N = 2$, and $\rho = 10$.

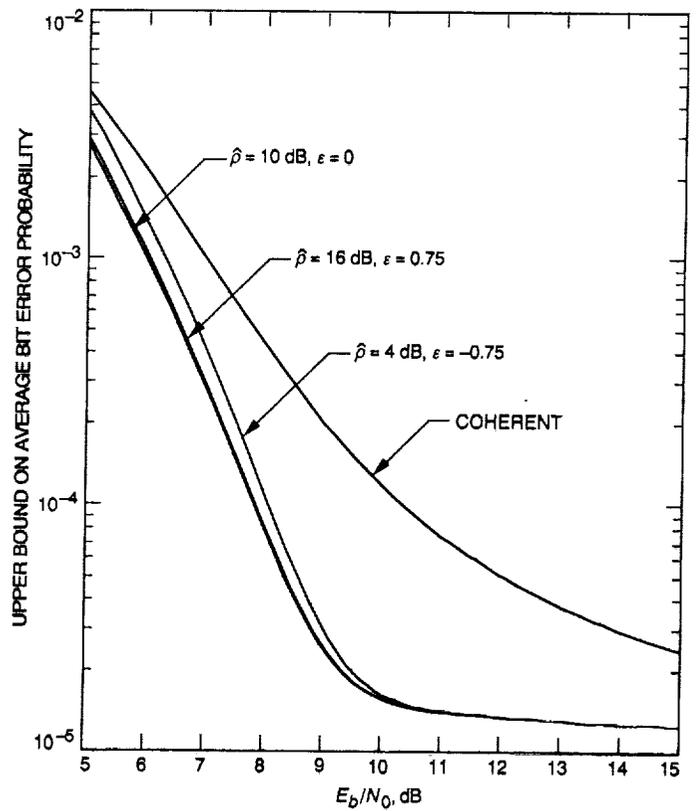


Fig. 7. Upper bound on average bit error probability versus E_b/N_0 in decibels for MLSE in the presence of loop SNR mismatch; $M = 2$, $N = 8$, and $\rho = 10$.

N 9 3 - 19432
128452
P-15

Parameter and Configuration Study of the DSS-13 Antenna Drives

W. Gawronski and J. A. Mellstrom
Ground Antennas and Facilities Engineering Section

The effects of different elevation and azimuth drive configurations on DSS-13 antenna performance are presented as well as a study of gearbox stiffness and motor inertia. Small motor inertia and rigid gearboxes would improve the pointing accuracy up to a certain limit. The limit is imposed by critical values of gearbox stiffness and motor inertia introduced in the article. The critical values depend on the lowest structural frequency of the rate-loop model. The tracking performance can be improved by raising gearbox stiffness to the critical stiffness and reducing motor inertia to the critical inertia. An azimuth drive configuration with four driven wheels was also investigated. For the four-wheel drive configuration in azimuth, the cross-coupling effects are reduced and wind disturbance rejection properties improved. Pointing is improved substantially in the cross-elevation but is relatively unaffected in the elevation direction. More significant improvements can be achieved through either structural redesign (stiffening the structure) or new control algorithms or control concepts, which would eliminate the effect of flexible deformations on the antenna pointing accuracy. Although the study is performed for the DSS-13 antenna, the results can be extended for other DSN antennas.

I. Introduction

This article investigates the DSS-13 antenna drives and their effect on antenna pointing accuracy. Each elevation and azimuth drive consists of a pair of motors and gearboxes. The size of a motor and a gearbox is determined from such criteria as static wind loads, which do not directly reflect pointing performance. The purpose of this study was to determine criteria for sizing motors and gearboxes so that the pointing accuracy is accounted for. For control system design purposes, the motor size is given in terms of motor inertia, while the gearbox size is given in terms of gearbox stiffness. Gearbox inertia is neglected since it is less than 10 percent of motor inertia when the

gear ratio is taken into account. Different locations of drives in azimuth and elevation are also investigated. One drive in elevation and two drives in azimuth are compared with two smaller drives in elevation and four smaller drives in azimuth, each at a different location. The effect of this drive configuration on antenna pointing accuracy is investigated.

II. Performance Criteria

Tracking performance and wind disturbance rejection are used to evaluate the pointing performance of motors and gearboxes in the elevation and azimuth drives. The

rate-loop bandwidth is used as a measure of tracking performance and rms pointing error due to wind gusts as a measure of wind disturbance rejection. For the purpose of this article, the rate-loop bandwidth is defined as a frequency range from zero up to the lowest lightly damped mode in the rate-loop transfer function (rate command to rate output). This definition is used for the PI controller design, and the lowest lightly damped mode determines the frequency range of the controller action. A lightly damped mode is detected as a resonant peak in the plot of magnitude of the transfer function (Fig. 1). The wider the open-loop bandwidth is, the better the closed-loop tracking performance is. The wind disturbance rejection properties are evaluated through simulations using the antenna model developed in [1] and the wind model described in [2].

III. Parameter Study

In this section, the effect of motor inertia and gearbox stiffness on antenna performance is investigated. It is obvious that a rigid drive would significantly improve a rigid antenna performance. For a flexible antenna, even a rigid drive cannot prevent its flexible deformations, thus the performance improvement through gearbox stiffening is limited. This is analyzed in detail below.

For the DSS-13 antenna performance evaluation (at a 60-deg elevation position), the model developed in [1] is used. The rate-loop model is shown in Fig. 2, where for clarity only the elevation drive is presented. The model consists of the antenna structure model (21 modes, up to 10 Hz, including two free-rotating modes), gearbox model, motor armature, and amplifiers. The elevation and azimuth drive configuration in the rate-loop model is shown in Fig. 3. The nominal motor inertia is $J_{mn} = 0.14 \text{ N m sec}^2$, (1.236 lb in sec^2), and the nominal gearbox stiffness is $k_{gn} = 1.65 \times 10^6 \text{ N m/rad}$ ($1.5 \times 10^7 \text{ lb in./rad}$).

The effect of the gearbox stiffness on the antenna performance is investigated by observing the change of the imaginary components of the rate loop-poles with respect to gearbox stiffness, see Fig. 4. The imaginary parts of the roots represent the structural natural frequencies. Natural frequencies of the structure and the gearbox are shown in this figure. The lowest structural frequency and the gearbox frequency define the bandwidth as shown in Fig. 4. The bandwidth grows with the gearbox stiffness, up to the critical value $k_{gc} = 1.1 \times 10^6 \text{ N m/rad}$ (10^7 lb in./rad). For

$$k_g > k_{gc} \quad (1)$$

the tracking performance remains unchanged. Thus, the critical stiffness k_{gc} defines the minimal stiffness of a gearbox, which assures reasonable tracking properties.

The bandwidth and the critical stiffness are also seen in the transfer function plots, Fig. 5. For $k_g < k_{gc}$, the gearbox resonant peak, which is smaller than the critical bandwidth, defines the bandwidth (Fig. 5a). For $k_g > k_{gc}$, the bandwidth changes insignificantly (Fig. 5b), since it is determined by the lowest natural frequency of the structure.

The wind disturbance rejection properties have been simulated for x -direction wind (along the elevation axis), and y -direction wind (horizontal direction orthogonal to the elevation axis). The results are summarized in Tables 1 and 2. The tables show that high gearbox stiffness improves wind rejection properties for y -direction wind, while the x -direction wind pointing remains almost unchanged.

The effect of motor inertia on antenna pointing performance is investigated in a similar fashion. The variations of the rate-loop poles due to motor inertia changes have been evaluated, and their imaginary parts are plotted in Fig. 6. The structural natural frequencies and the gearbox frequency are distinguished in this plot. The lowest natural frequency of the structure and the gearbox frequency determine the bandwidth. For

$$J_m > J_{mc} \quad (2)$$

the bandwidth is constant, and decreases for $J_m > J_{mc}$, deteriorating the antenna tracking properties. Thus, $J_{mc} = 50 \text{ lb in sec}^2$ is the critical value of inertia. The phenomenon can be observed in the transfer function plots (elevation rate command to elevation rate), Fig. 7, where for small inertia (small when compared with the critical one) the bandwidth is constant and for large inertia the bandwidth narrows.

Tables 1 and 2 summarize wind disturbance rejection properties. They show that the properties do not improve with motor inertia decrease below the critical value.

IV. Configuration Study

The existing drive configuration of the DSS-13 antenna is shown in Fig. 3. It consists of one elevation and two azimuth drives. A new configuration is compared. The number of drives in this configuration is doubled, and they are mounted at different structural locations. Motors are

sized so that their total power is the same as in the original configuration.

Due to the high stiffness of the bullgear, two elevation drives at different locations of the bullgear have the same effect as two drives at the same location. Therefore, two drives are equivalent to one drive with a properly sized motor (the motor inertia of the two-drive configuration is 38 percent of the motor inertia of the one-drive configuration). Hence the two-elevation drive case reduces to the one-drive parameter study presented previously. The transfer function plots in Fig. 7 compare one- and two-elevation drive cases. They show that the bandwidth in azimuth and elevation remains the same. Thus, no improvement in tracking accuracy is observed. Also, simulations show no improvement in the x -direction wind disturbance rejection and moderate improvement in y -direction wind disturbance rejection (assuming rigid enough gearboxes).

Since the stiffness of the alidade is comparable with the stiffness of gearboxes, the problem of four-azimuth drives cannot be reduced to an equivalent two-drive problem. In the four-azimuth-drive configuration, each drive is mounted on a separate azimuth wheel. The same gearbox stiffness is assumed, and the motor inertia is 2.6 times smaller than the motor inertia of the two-drive case. The tracking performance is evaluated through bandwidth comparison of two- and four-azimuth drives (Fig. 8), and through step response simulations (Fig. 9). In Fig. 8(a) the bandwidth in azimuth is slightly larger, and in elevation it remains the same in the four-drive case, while the cross-transfer function (from elevation to azimuth rate and from azimuth to elevation rate) shows significant change. It is confirmed by the closed-loop unit step responses. The responses to an azimuth step command differ slightly for two- and four-azimuth drives, and the responses to an elevation step command overlap in both cases, while cross-coupling responses show significant differences between the four- and two-drive case.

Wind disturbance rejection for the two- and four-azimuth-drive case is compared in Fig. 10 and Tables 3

and 4. The tables show improvement in x -direction wind rejection properties in the four-drive case, but additional stiffening of drives does not improve the wind disturbance rejection properties.

V. Conclusions

The article has defined criteria for drive comparison purposes and determines conditions imposed on gearbox stiffness and motor inertia so that the tracking errors are minimized and wind disturbance rejection properties are improved. It showed the critical values of gearbox stiffness and motor inertia limit tracking performance improvement. The gearbox stiffness should be larger (but not necessarily much larger) than the critical stiffness, and the motor inertia should be smaller (but not necessarily much smaller) than the critical inertia in order to preserve tracking accuracy. The existing (nominal) parameters of the DSS-13 antenna satisfy these demands. An overdesigned drive is a drive with a gearbox stiffness much larger than the critical one and/or motor inertia much smaller than the critical one. Overdesigned drives do not significantly improve the tracking performance, although an overdesigned gearbox improves wind disturbance rejection. Also, the four-azimuth-drive configuration does not improve the tracking performance (bandwidth remains almost unchanged), but improves the cross-dynamic properties and wind disturbance rejection properties for winds from y -direction.

Improvements due to stiffening gearboxes to downsizing drive motors, and to multiple drives are non-negligible, but not dramatic. Thus, for moderate improvement of performance it is advised to stiffen the gearboxes and use four-azimuth drives. Significant improvement may be achieved only through more innovative approaches, such as antenna structure redesign (more rigid), application of a new control algorithm (with vibration suppression properties), or implementation of either new or additional sensors/actuators (e.g., active truss members for structural vibration damping).

References

- [1] W. Gawronski and J. A. Mellstrom, "Modeling and Simulations of the DSS-13 Antenna Control System," *TDA Progress Report 42-106*, vol. April-June 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 205-248, August 15, 1991.
- [2] W. Gawronski, B. Bienkiewicz, and R. E. Hill, "Pointing-Error Simulations of the DSS-13 Antenna Due to Wind Disturbances," *TDA Progress Report 42-108*, vol. October-December 1991, Jet Propulsion Laboratory, Pasadena, California, pp. 109-134, February 15, 1992.

Table 1. Pointing errors due to x-direction wind.

Drive parameters	k_{gn}, J_{mn}	$10k_{gn}, J_{mn}$	$k_{gn}, 0.5J_{mn}$
Elevation pointing error, mdeg	2.81	2.30	2.89
Cross-elevation pointing error, mdeg	1.64	2.01	1.69
X-band loss, dB	0.03	0.03	0.03
Ka-band loss, dB	0.44	0.39	0.46

Table 2. Pointing errors due to y-direction wind.

Drive parameters	k_{gn}, J_{mn}	$10k_{gn}, J_{mn}$	$k_{gn}, 0.5J_{mn}$
Elevation pointing error, mdeg	3.77	2.09	4.10
Cross-elevation pointing error, mdeg	0.55	0.32	0.77
X-band loss, dB	0.04	0.01	0.05
Ka-band loss, dB	0.60	0.19	0.72

Table 3. Pointing errors due to x-direction wind disturbances for two- and four-azimuth drives.

Drive parameters	2AZ, k_{gn}	4AZ, k_{gn}	2AZ, $10k_{gn}$	4AZ, $10k_{gn}$
Elevation pointing error, mdeg	2.81	2.43	2.51	2.40
Cross-elevation pointing error, mdeg	1.64	0.43	2.08	0.36
X-band loss, dB	0.03	0.02	0.03	0.02
Ka-band loss, dB	0.44	0.25	0.44	0.24

Table 4. Pointing errors due to y-direction wind disturbances for two- and four-azimuth drives.

Drive parameters	2AZ, k_{gn}	4AZ, k_{gn}	2AZ, $10k_{gn}$	4AZ, $10k_{gn}$
Elevation pointing error, mdeg	3.77	3.77	3.51	3.92
Cross-elevation pointing error, mdeg	0.55	0.53	0.34	0.50
X-band loss, dB	0.04	0.04	0.04	0.04
Ka-band loss, dB	0.60	0.60	0.52	0.65

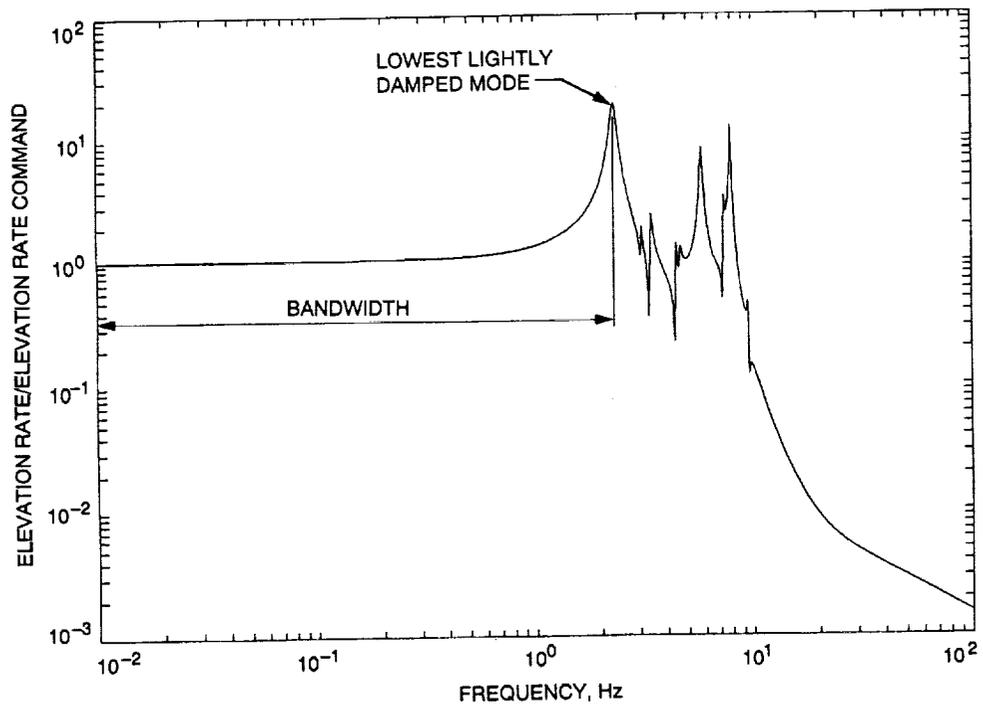


Fig. 1. Bandwidth of the antenna.

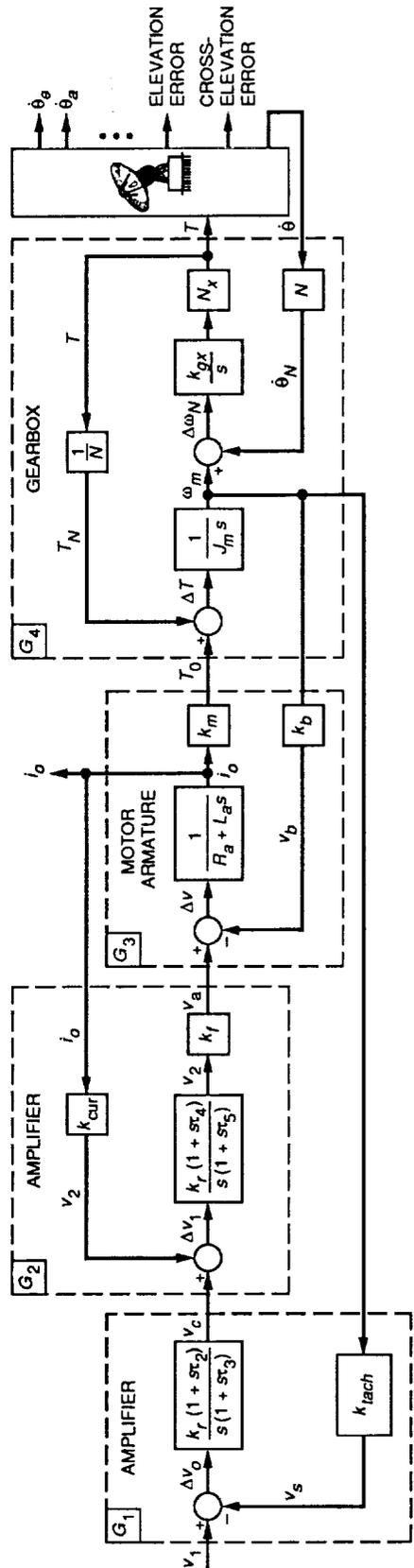


Fig. 2. Rate-loop model of the antenna with the elevation drive.

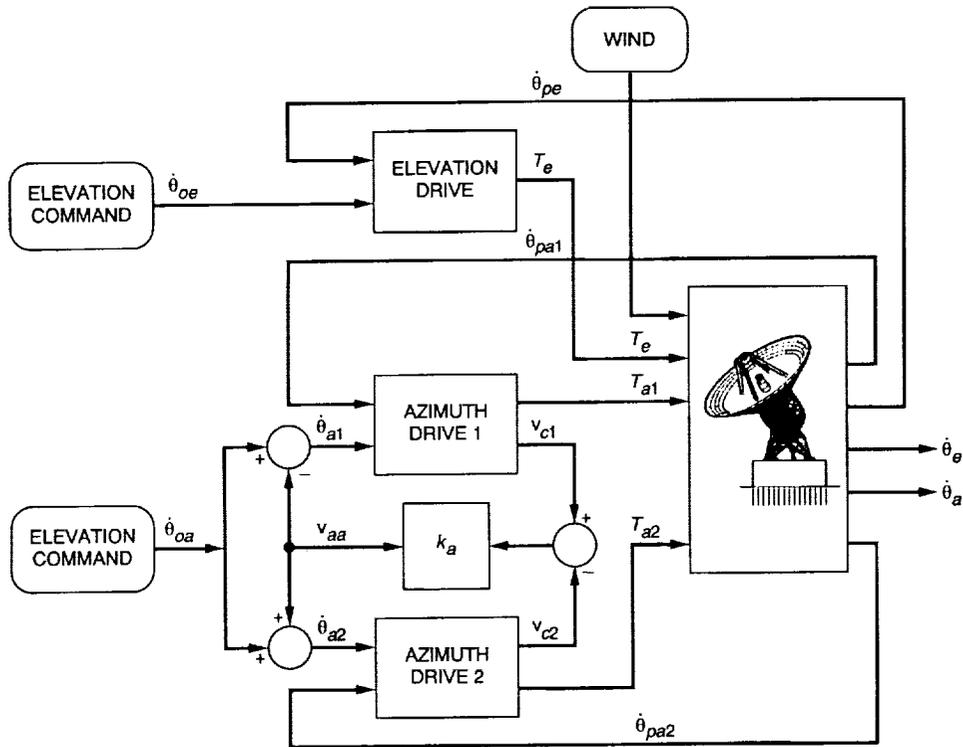


Fig. 3. Rate-loop model of the antenna.

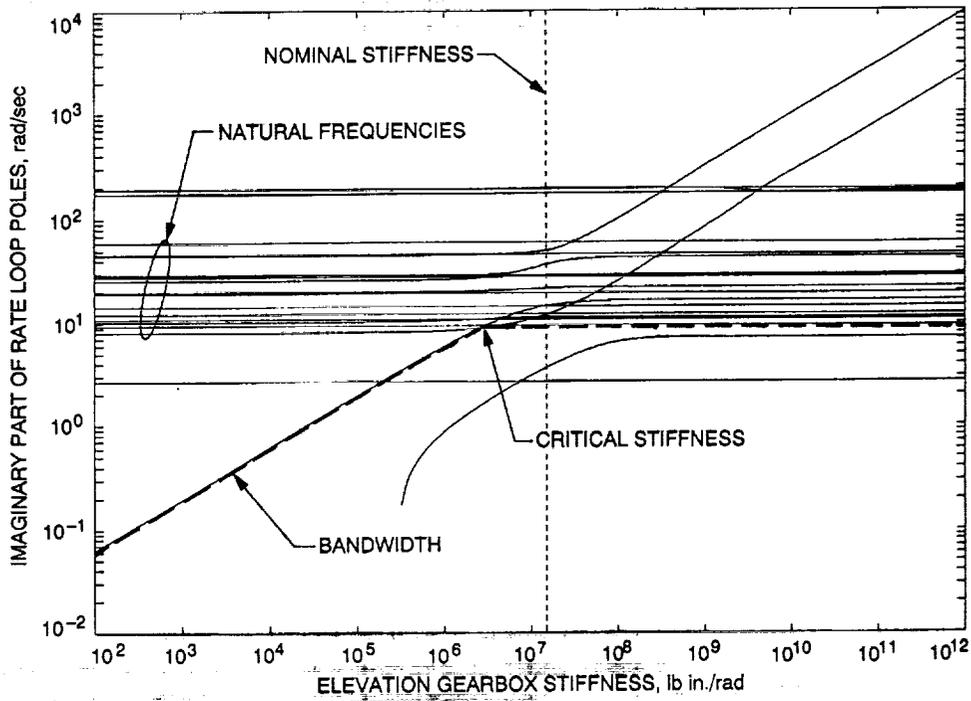


Fig. 4. Imaginary parts of the rate loop poles versus elevation gearbox stiffness.

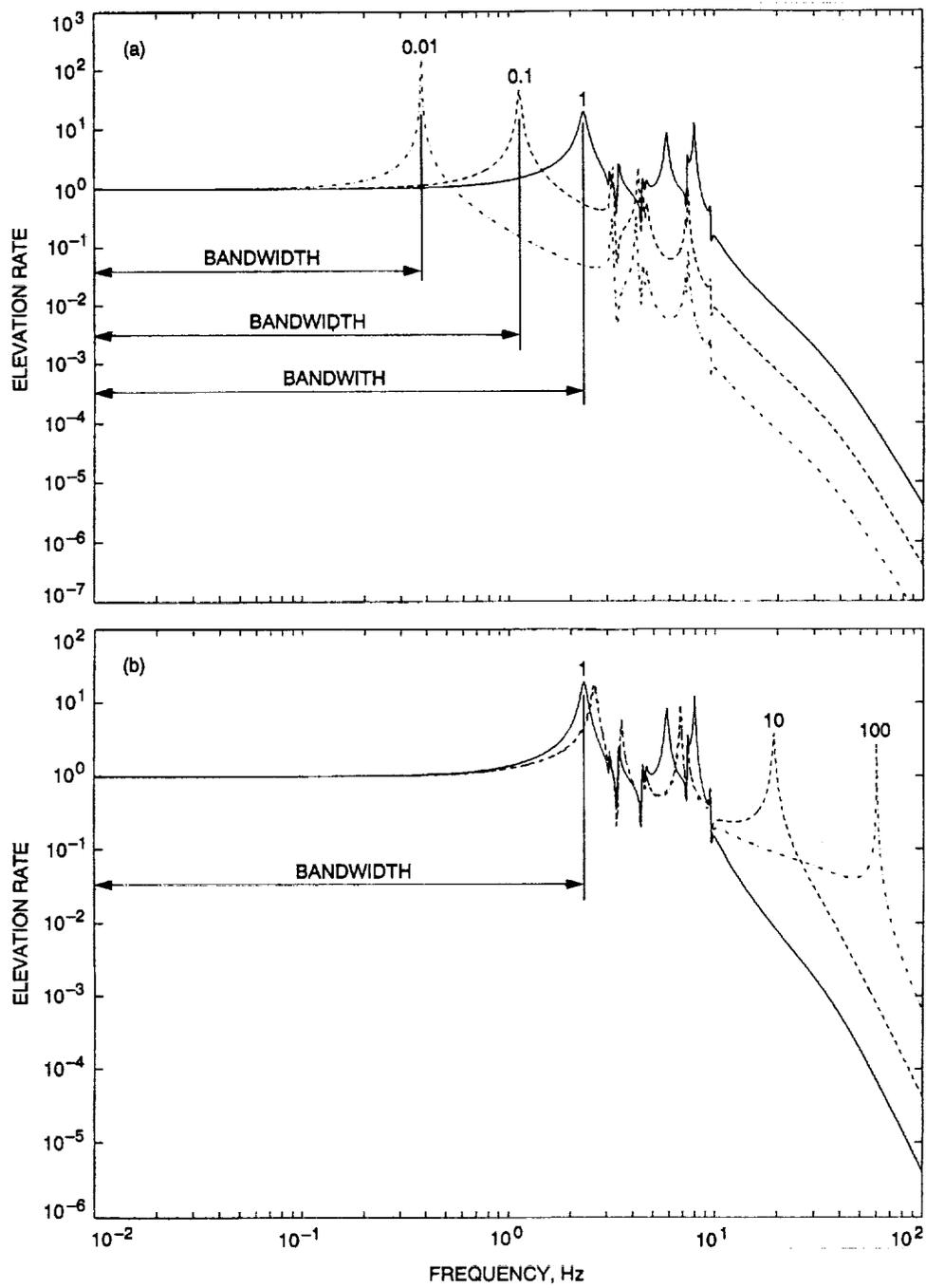


Fig. 5. Transfer functions: elevation rate output to elevation rate command.

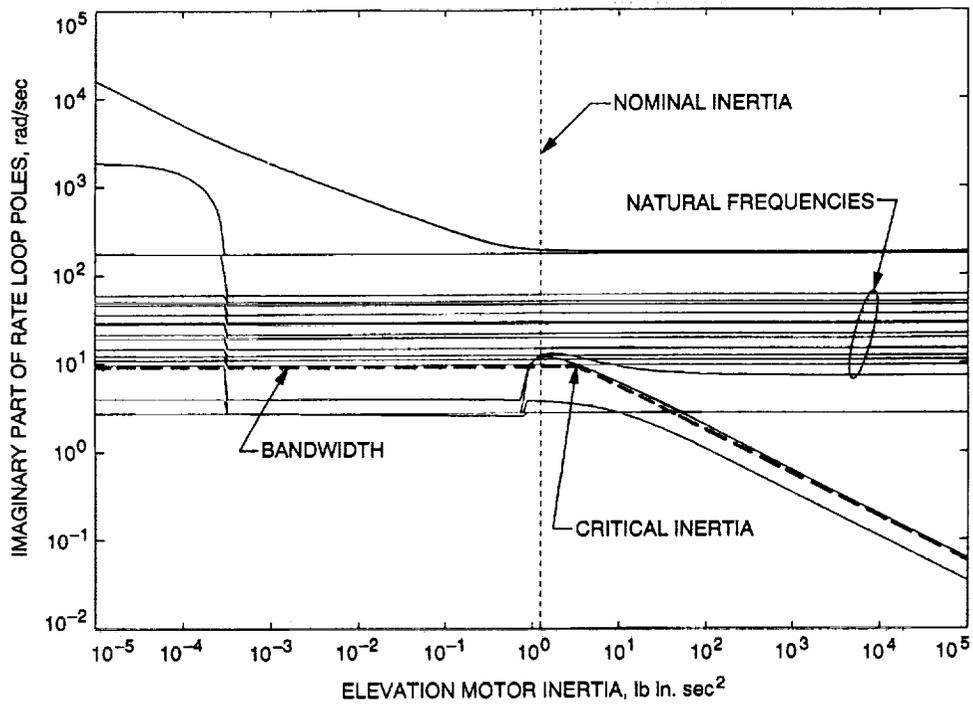


Fig. 6. Imaginary parts of the rate-loop poles versus elevation motor inertia.

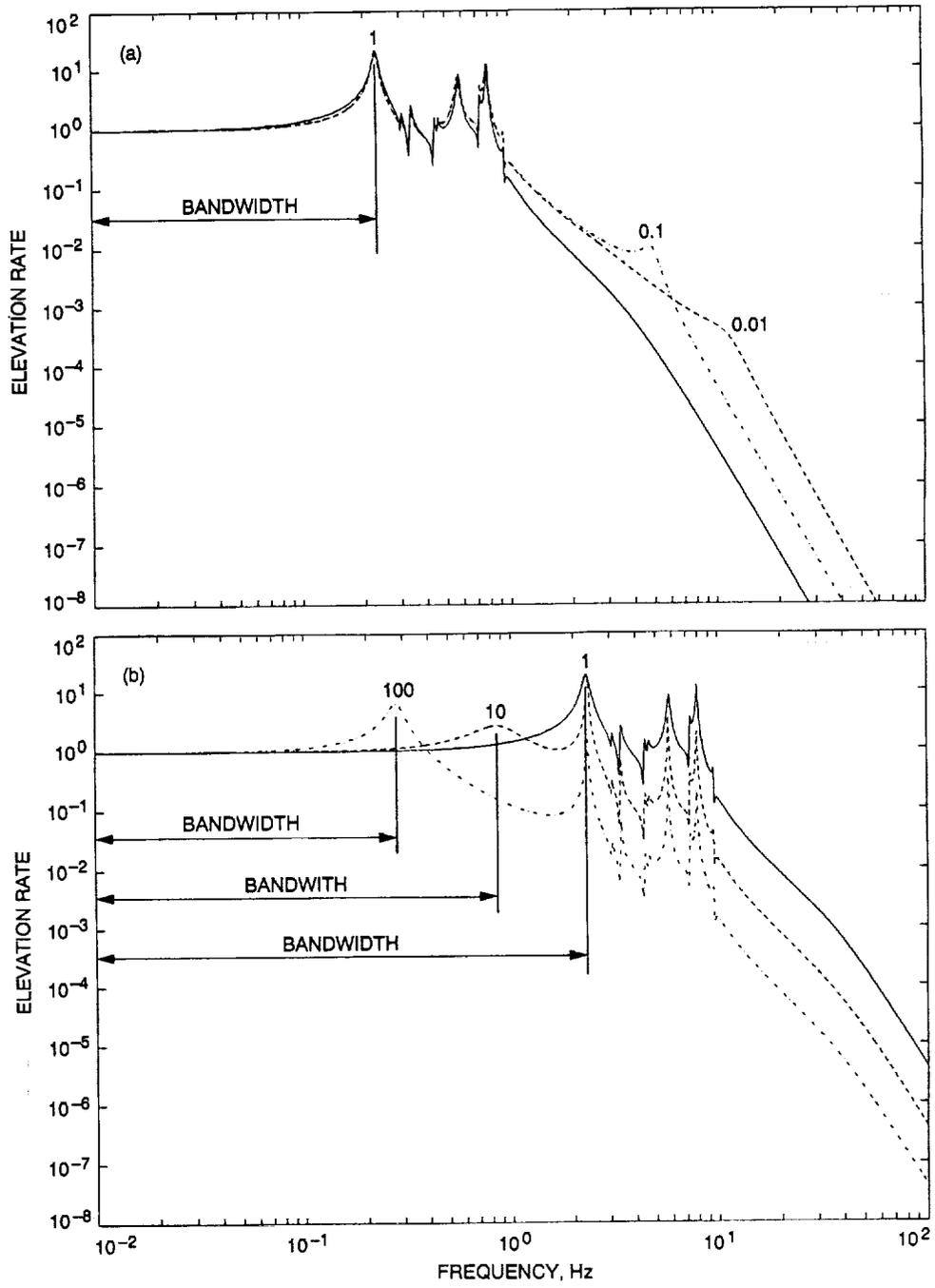


Fig. 7. Transfer functions: elevation rate output to elevation rate command.

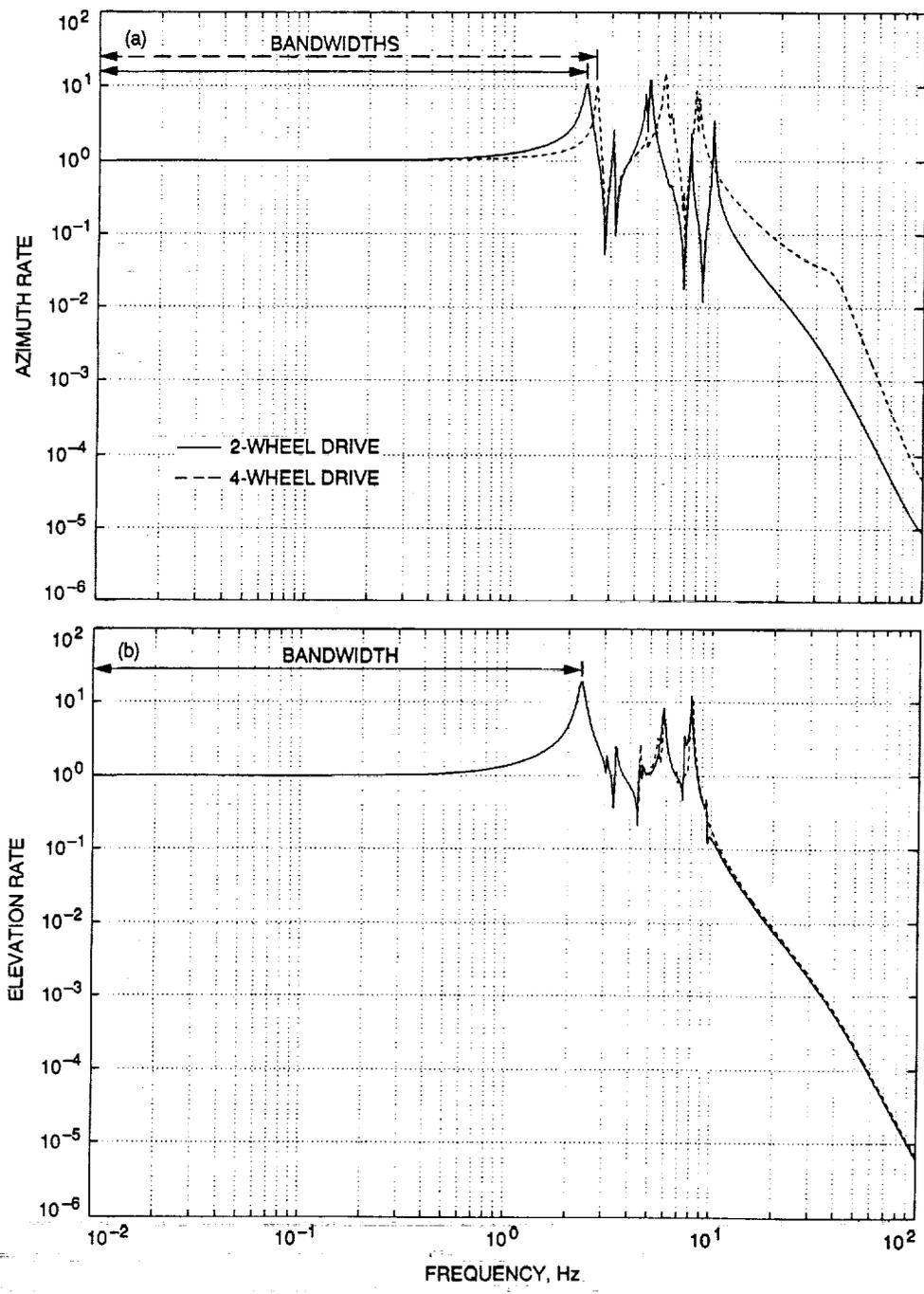


Fig. 8. Bandwidth for two- and four-wheel azimuth drives: (a) azimuth rate to azimuth rate command and (b) elevation rate to elevation rate command.

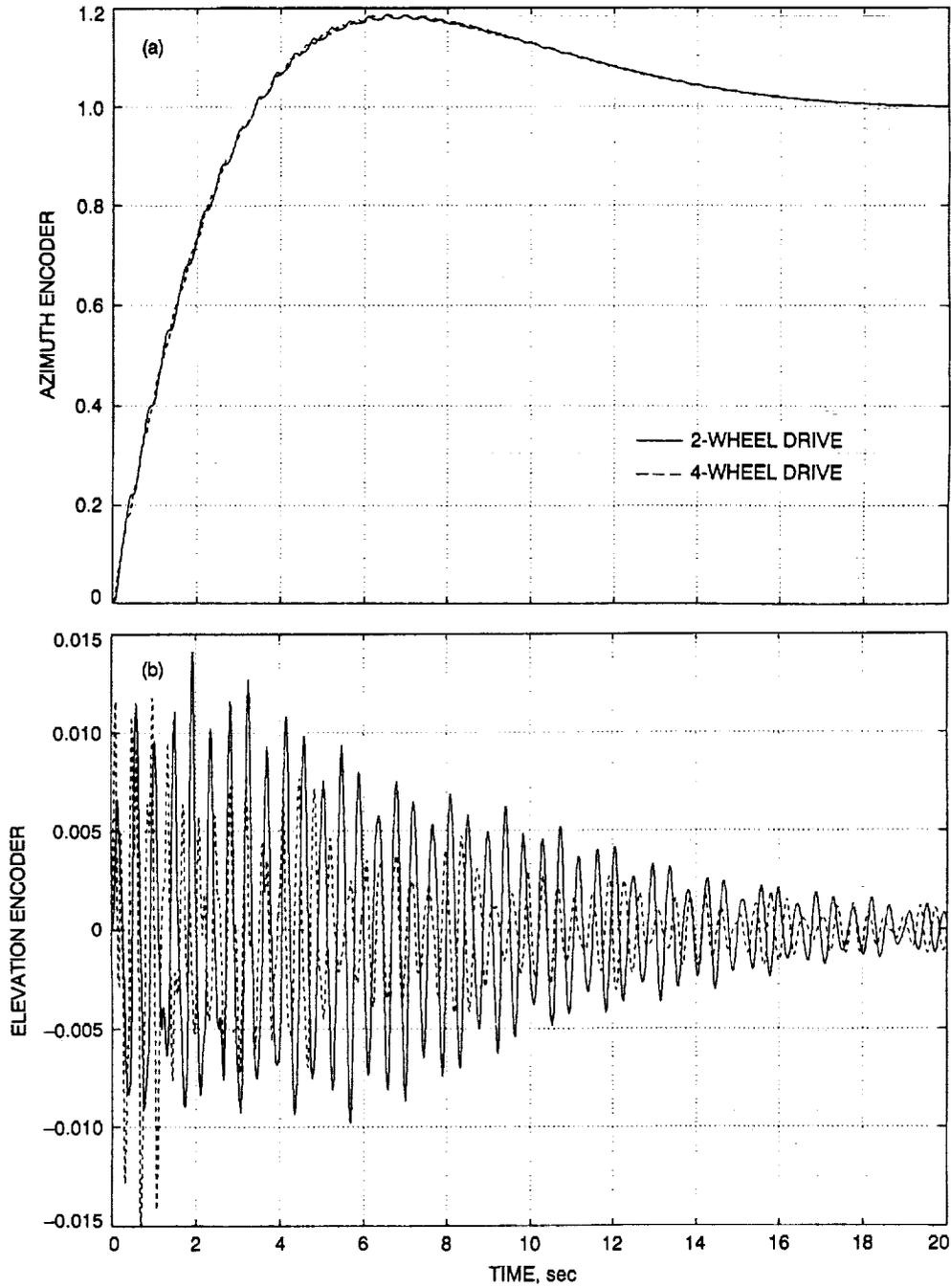


Fig. 9. Step responses for two- and four-wheel azimuth drives: (a) azimuth encoder to azimuth command; (b) elevation encoder to azimuth command; (c) elevation encoder to elevation command; and (d) azimuth encoder to elevation command.

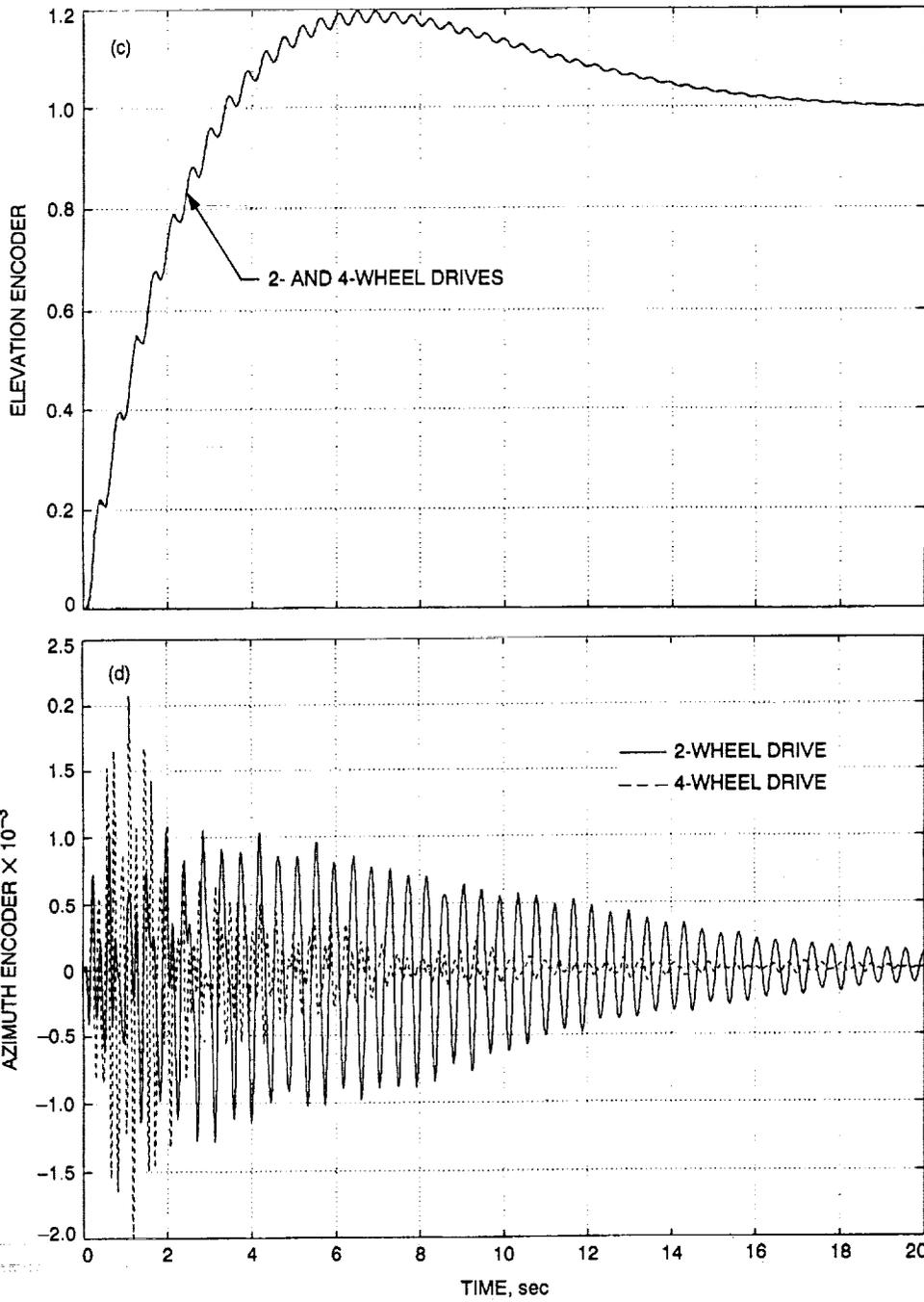


Fig. 9 (contd).

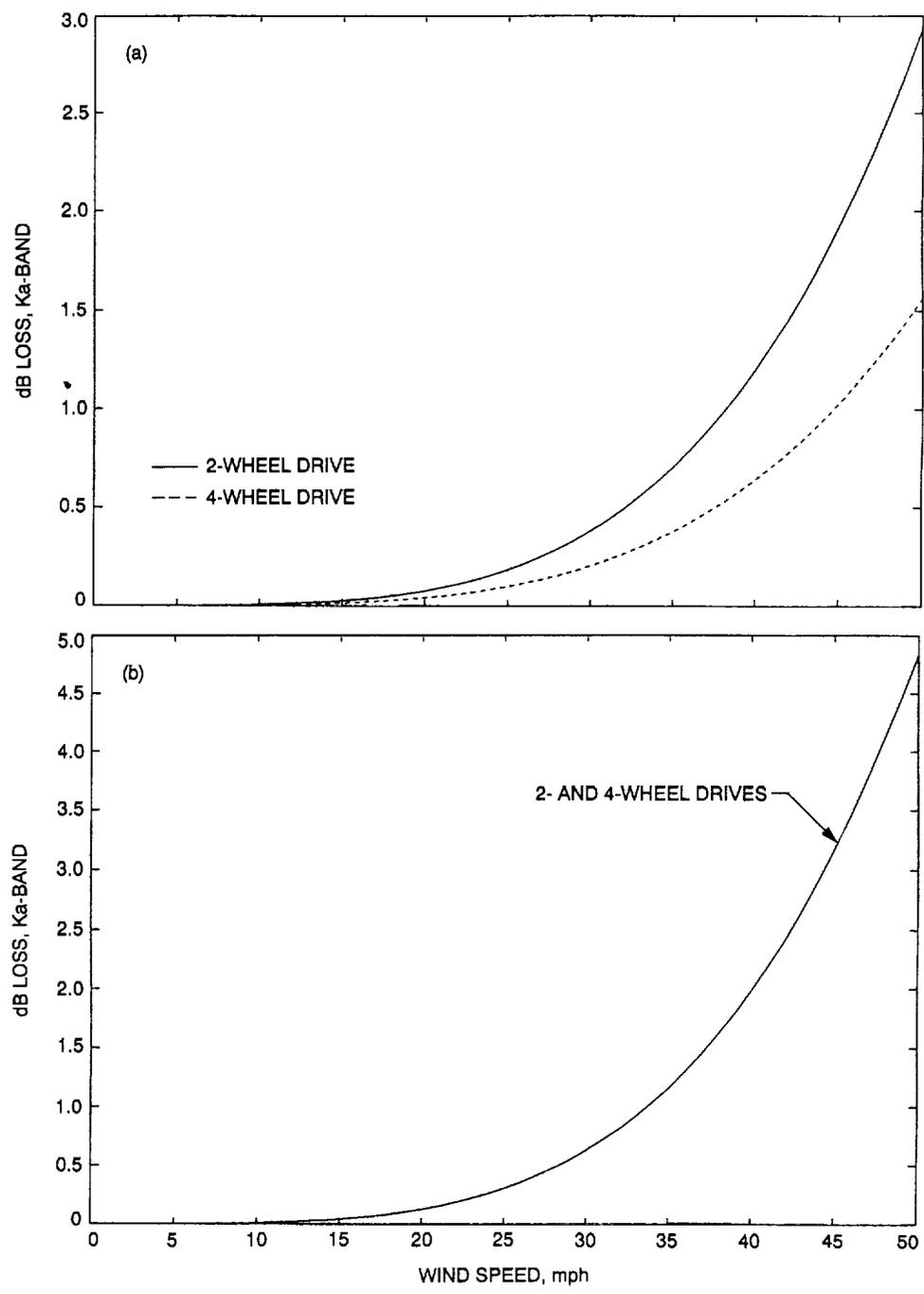


Fig. 10. Ka-band decibel loss due to wind gusts: (a) x-direction wind and (b) y-direction wind.

520-32
 123456789
 10111213141516171819202122232425262728293031323334353637383940414243444546474849505152535455565758596061626364656667686970717273747576777879808182838485868788899091929394959697989900
 P-11

SETI Low-Frequency Feed Design Study for DSS 24

P. H. Stanton and P. R. Lee
 Ground Antennas and Facilities Engineering Section

The Search for Extraterrestrial Intelligence Sky Survey project requires operation from 1 to 10 GHz on the beam waveguide (BWG) antenna DSS 24. The BWG reflectors are undersized in the 1- to 3.02-GHz range, resulting in poor performance. Horn designs and a method for implementing 1- to 3.02-GHz operation on DSS 24 are presented. A combination of a horn and a shaped feed reflector placed above the main reflector is suggested. The horn and feed reflector could be hidden in the RF shadow of the subreflector and struts. Results from computer analysis of this design indicate that adequate performance could be achieved.

I. Introduction

The DSS 24 34-m beam waveguide (BWG) antenna will be used for the Search for Extraterrestrial Intelligence (SETI) Sky Survey project over a frequency range of 1 to 10 GHz. The antenna's BWG reflectors are undersized for low noise operation at the low end of this frequency range. Therefore, an alternate method of feeding the antenna from 1 to 3.02 GHz is presented in this article. A corrugated, 29-dB-gain horn located at the Cassegrainian focus would give acceptable RF performance but would be physically large and would block the normal BWG transmission path. To reduce the size of the feed horn and facilitate swift clearing of the BWG path for normal DSN operation, a smaller feed horn in combination with a movable shaped reflector could be used instead. This feed reflector would sit over the BWG aperture in DSS 24's main reflector during SETI operations (see Fig. 1(a)) and would be moved into a storage position during DSN operations.

The SETI horn-reflector combination could be hidden in the RF shadow of the subreflector and struts as shown in Fig. 1(b). Major RF system requirements¹ are listed in Table 1.

II. Horn Design and Analysis

Two small-aperture, corrugated horns were used to cover the lower frequency range of the SETI sky survey (horns number one and number two would operate over the ranges of 1 to 1.73 GHz and 1.61 to 3.02 GHz, respectively). The aperture diameter of each of these horns was constrained to three wavelengths, at their lowest operating frequency, in order to help limit its weight and cost. The

¹ G. A. Zimmerman, *Search for Extra-Terrestrial Intelligence Microwave Observing Project Sky Survey Element*, 1720-4100 (internal document), Jet Propulsion Laboratory, Pasadena, California, November 12, 1991.

output flare angles of the horns were adjusted along with the shape of the feed reflector to empirically arrive at an acceptable performance across the two frequency ranges.

Once the frequency range, aperture size, and output flare angle were established, the detailed wideband horn design followed the procedures outlined in [1]. The corrugation profile of the 1- to 1.73-GHz horn is shown in Fig. 2, with its various sections labeled in accordance with [1] (mode converter, frequency transition, angle transition, and output flare).

The corrugated feedhorns were analyzed with a JPL computer program² that uses modal field-matching techniques to determine the transverse electric (TE_{mn}) and transverse magnetic (TM_{mn}) scattering matrix of the horns [2]. From the scattering matrix, the return loss and aperture fields were known. Using the radiation integral, the radiation patterns of the horns were calculated from the aperture fields. Table 2 summarizes the performance of the horns at several frequencies.

III. Feed Reflector Design and Analysis

The shaped feed reflector was designed to optimize system performance at the lowest frequencies, where system spillover was expected to be the highest and to result in higher noise temperature and lower gain. The center of the feed reflector was chosen so that the feed reflector would be located a little above the main reflector. To approximate the location of the first geometric optics (GO) focal point of the feed reflector, a 25-dB horn pattern was used to illuminate the subreflector. This pattern was moved along the axis of the subreflector, and the overall gain of the system was calculated using JPL computer programs³ that employ physical optics (PO) techniques. At the location resulting in the best gain, the far-field phase center of the horn pattern was 437 cm above F1, the focal point of the main reflector and subreflector system [see Fig. 1(a)]. The first GO focal point of the feed reflector was placed 445 cm from the reflector. The second GO focal point was placed in the desired location of the feedhorn far-field phase center. This location was chosen so that the horns would be close to the feed reflector without blocking reflected radiation. The resulting shape was an ellipsoid.

² D. J. Hoppe, *Scattering Matrix Program for Ring-Loaded Circular Waveguide Junctions* (internal document), Jet Propulsion Laboratory, Pasadena, California, August 3, 1987.

³ R. E. Hodges and W. A. Imbriale, *Computer Program POMESH for Diffraction Analysis of Reflector Antennas* (internal document), Jet Propulsion Laboratory, Pasadena, California, February 1992.

The reflector parameters, horn location, and horn flare angle were varied until the performance was acceptable in the 1- to 1.73-GHz range. The performance was then evaluated in the 1.61- to 3.02-GHz range. The horn location and flare angle were varied to achieve acceptable performance. Other reflector parameters were also tried, but the performance was not improved. Figure 3 shows the final reflector design.

The main reflector and subreflector system was designed to work optimally with a 29-dB pattern that has its near-field phase center at F1 and using an observation distance to the subreflector. How the system performs with the GO focal point of the feed reflector so far above F1 was investigated. The phase centers (PC's) of the 1-GHz and 3.02-GHz patterns generated from the horn and feed reflector combination were calculated in the far field and at various distances (R equals observation radius) in the near field (see Fig. 4). As the observation distance moved further into the near field, the phase center moved along the z -axis in a negative direction. Defocus curves of antenna system gain versus feed reflector location were generated by moving the patterns along the axis of the antenna system (see Fig. 5). At the high-frequency limit, the far-field phase center of a focused system is located at the focal point of the ellipsoidal reflector, which is 508 cm above the reflector.

At 1 GHz, the far-field phase center is 208 cm above the feed reflector. The reflector is only about 10 wavelengths in diameter at 1 GHz, so the system will not focus well in any case, but it appears to be somewhat in focus. With the feed reflector at the design location and using an observation distance to the subreflector, the near-field phase center is 7 cm below F1. Using the best gain location from Fig. 5(a), which yields only a 0.05-dB increase in gain, the near-field phase center is 69 cm below F1. Since a wavelength is 30 cm long, the phase centers are reasonably close to F1. At this long wavelength, the system is fairly insensitive to movement of the feed reflector along the main reflector axis.

At 3.02 GHz, the horn and feed reflector system is out of focus, causing the far-field phase center of the reflected pattern to be 155 cm below the reflector. With the feed reflector at the design location, and using an observation distance to the subreflector, the near-field phase center is 320 cm below F1. With a wavelength of 9.9 cm at 3.02 GHz, the phase center is very far from F1. Performance at this frequency could be improved by redesigning the feed reflector to move the near-field phase center close to F1. Simply moving the feed reflector does not improve performance, as shown in the defocus curve of Fig. 5(b).

In an attempt to focus the system some and increase the overall gain, the subreflector was moved along its axis. By moving the subreflector 1.6 cm towards the main reflector, the gain was increased by 0.16 dB.

IV. Results

Table 3 summarizes the antenna system performance as calculated with the PO programs. In the 1- to 1.73-GHz band, the system performs well at one horn location. In the 1.61- to 3.02-GHz band, however, the horn must be moved to four different locations to achieve acceptable performance.

A study of the BWG at 1 GHz was performed using $\cos^n(\Theta)$ patterns as the horn input. The best gain-to-noise temperature ratio (G/T) achieved was 28.24 dB, with a gain of 47.74 dB and a noise temperature of 83.16 K. Compared to the results from horn number one plus the shaped reflector, the noise temperature increases 56.69 K, the gain decreases 2.02 dB, and the G/T decreases 7.30 dB.

V. Conclusion

DSS 24, operating in its normal configuration with the BWG, performs unacceptably at the low frequencies in the SETI 1- to 10-GHz range. By putting a horn and shaped reflector above the main reflector, the BWG could be bypassed, and acceptable performance could be achieved. The horn and shaped reflector would have minimal impact on the other DSS 24 operations since the horn and reflector would be placed in the RF shadow of the subreflector and struts.

The results presented in this article indicate that the suggested configuration would meet SETI system requirements with the two exceptions of noise temperature and beam efficiency. The noise temperature is slightly high in the lower frequencies of both the 1- to 1.73-GHz band and the 1.61- to 3.02-GHz band; however, the G/T is higher than the target G/T over the entire frequency range. The beam efficiency is lower than the system requirements over most of the 1- to 3.02-GHz range. The minimum beam efficiency is 81 percent at 2.21 GHz, and the system requirement is 90 percent over the entire range.

References

- [1] B. M. Thomas, G. L. James, and K. L. Greene, "Design of Wide-Band Corrugated Conical Horns for Cassegrain Antennas," *IEEE Transactions on Antennas and Propagation*, vol. AP-34, no. 6, pp. 750-757, June 1986.
- [2] G. L. James and B. M. Thomas, "TE₁₁ to HE₁₁ Cylindrical Waveguide Mode Converters Using Ring-Loaded Slots," *IEEE Transactions on Microwave Theory and Techniques*, vol. MTT-30, no. 3, pp. 278-285, March 1982.

Table 1. System requirements.

Parameter	Required value
Noise temperature	≤ 25 K
Polarization	RCP ^a and LCP ^b (simultaneous)
Instantaneous bandwidth	≥ 360 MHz
Aperture efficiency	≥ 65 percent
Beam efficiency	≥ 90 percent

^a Right-circular polarization.

^b Left-circular polarization.

Table 2. Horn performance.

Frequency, GHz	Gain, dB	Return loss, dB (TE_{11} mode)	Maximum cross-polarization, dB
1- to 1.73-GHz horn			
1.00	17.0	-33	-33.1
1.15	17.2	-41	-30.6
1.32	18.0	-51	-42.7
1.51	18.6	-49	-37.0
3.02	19.5	-47	-37.3
1.61- to 3.02-GHz horn			
1.61	17.1	-65	-29.8
1.88	18.2	-46	-35.6
2.21	19.0	-50	-47.0
2.58	19.7	-47	-37.8
3.02	20.6	-48	-41.7

Table 3. System performance.

Frequency, GHz	Noise Temperature ^a , K	Gain, dB	G/T , dB/K	Target G/T ^b , dB/K	Aperture efficiency	Beam efficiency
1.0- to 1.73-GHz horn						
1.00	26.47	49.77	35.54	35.18	0.747	0.861
1.15	24.68	50.93	37.01	36.40	0.739	0.871
1.32	22.06	52.08	38.65	37.59	0.731	0.915
1.51	21.10	53.11	39.87	38.76	0.708	0.935
1.73	20.55	54.19	41.06	39.94	0.691	0.945
1.61- to 3.02-GHz horn						
1.61	26.52	53.60	39.37	39.32	0.698	0.814
1.88	24.68	54.92	40.99	40.66	0.692	0.854
2.21	24.02	56.47	42.66	40.07	0.716	0.809
2.58	23.46	57.81	44.10	43.41	0.715	0.830
3.02	23.03	59.10	45.47	44.78	0.703	0.867

^a The noise temperature includes contributions from the sky and atmosphere, the reflectors, the horn, and the low-noise amplifier (LNA) assembly. Sky and atmosphere and LNA assembly values from: JPL-ARC Front-End Design Team, *NASA SETI Common Radio Frequency System Design Team Report* (NASA internal report), Appendix D, p. 6, NASA, Washington, DC, August 1, 1991.

^b The target G/T is for 65-percent aperture efficiency and 25-K noise temperature.

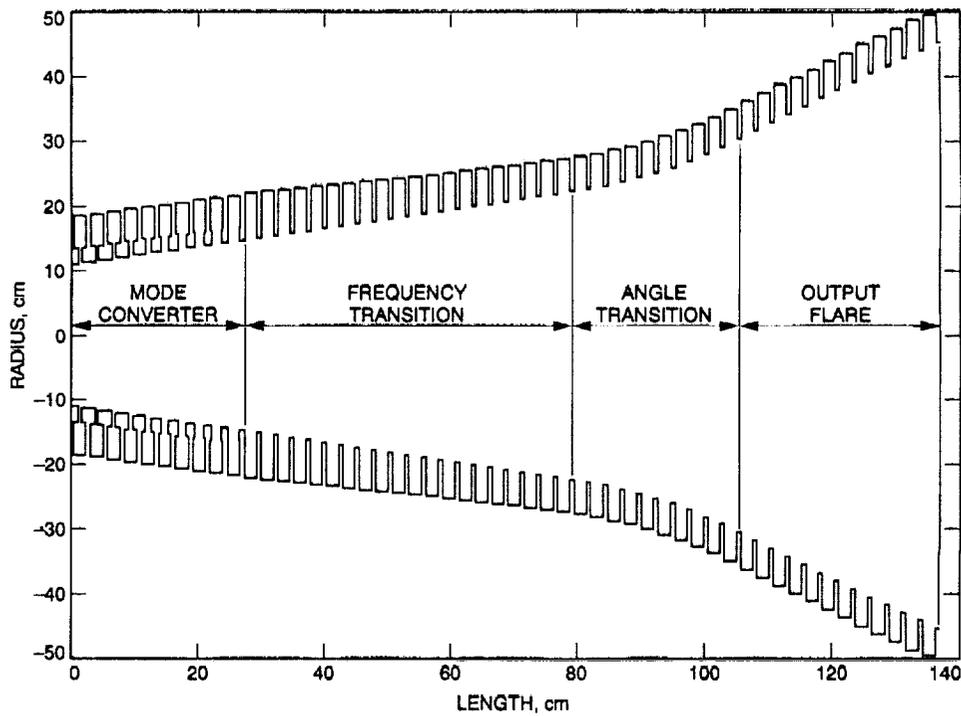


Fig. 2. Profile of 1- to 1.73-GHz horn.

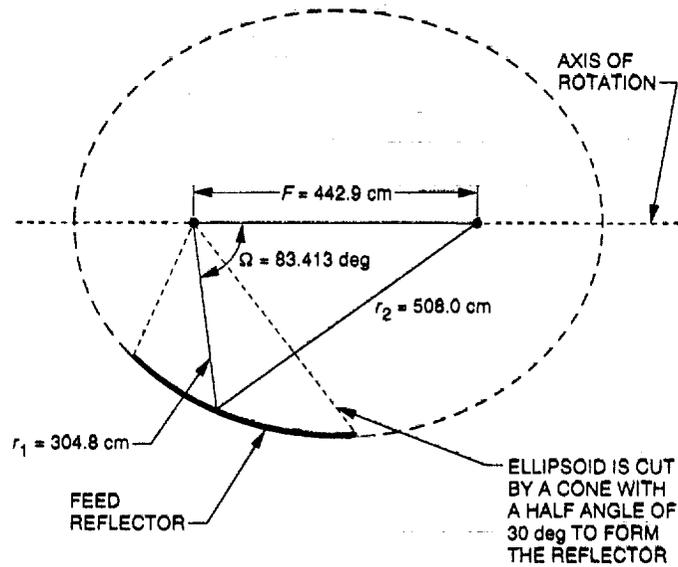


Fig. 3. Feed reflector.

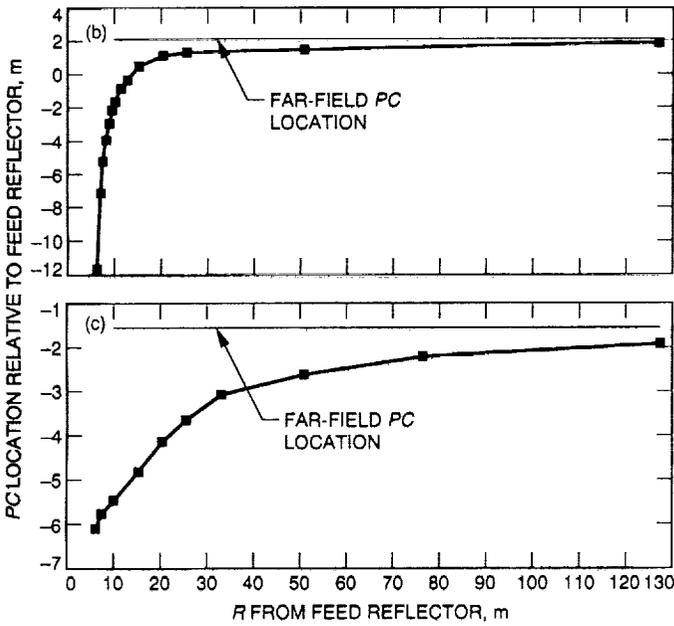
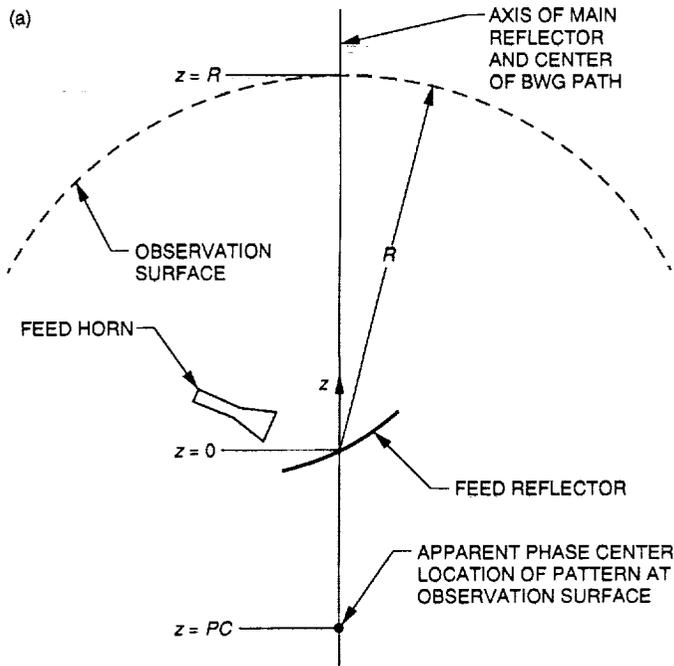


Fig. 4. Phase center location: (a) for horn and feed reflector pattern; (b) versus observation radius at 1 GHz; and (c) versus observation radius at 3.02 GHz.

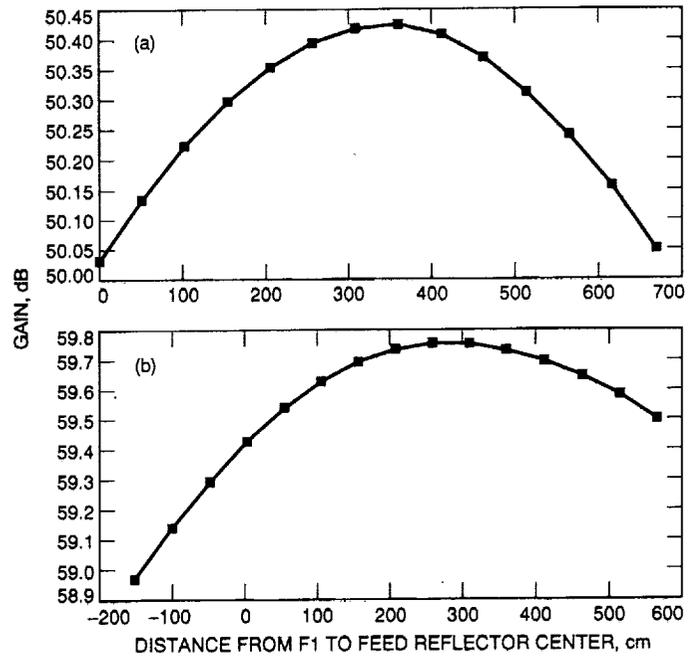


Fig. 5. Gain versus feed reflector location at: (a) 1 GHz and (b) 3.02 GHz.

N 9 3 - 1 9 4 3 4

S21-32

128454

R-10

Feed-Forward Control Upgrade of the Deep Space Network Antennas

W. Gawronski

Ground Antennas and Facilities Engineering Section

In order to improve the accuracy of high-rate tracking of NASA's DSN antennas, the position-loop controller has been upgraded with a feed-forward loop. Conditions for perfect and approximate tracking with the feed-forward loop are presented. The feed-forward loop improves tracking performance and preserves wind disturbance rejection properties of the previous closed-loop system.

Pointing accuracy of a proportional and integral (PI) control system for the DSN antennas [1] is satisfactory for slow-tracking antennas but significantly deteriorates when tracking fast-moving objects. In order to improve the tracking accuracy in the latter case, a PI control system has been augmented with a feed-forward loop, as shown with the block diagram in Fig. 1. In this diagram, G_p , G_c , G_f , and G_w denote transfer functions of the antenna's rate loop, PI controller, feed-forward gain, and wind disturbance, respectively; and r is a command, y is output (elevation and azimuth angles), e is tracking error in azimuth and elevation, u is plant input, and w is wind disturbance. The plant transfer function $G_p(\omega)$ is a 2×2 matrix, with elevation and azimuth rates as inputs and elevation and azimuth angles as outputs.

In order to analyze the impact of the feed-forward gain on the closed-loop system performance, the transfer function from the command r and wind disturbance w to the tracking error e was derived. From Fig. 1, one obtains

$$e = r - y \quad (1a)$$

$$y = G_p u + G_w w \quad (1b)$$

$$u = G_f r + G_c e \quad (1c)$$

Assuming $I + G_p G_c$ to be nonsingular and denoting that $G_o = (I + G_p G_c)^{-1}$, from Eqs. (1a), (1b), and (1c), one obtains

$$e = G_o(I - G_p G_f)r - G_o G_w w \quad (2)$$

From the above equation one obtains perfect tracking (i.e., $e = 0$) in the absence of wind disturbances for the feed-forward gain G_f such that

$$G_p(\omega)G_f(\omega) = I \quad (3)$$

In the case of the DSN antennas, the condition (3) can be satisfied in a certain frequency range only. Simply by inspection of the magnitudes of the plant transfer function in Fig. 2(a-d), one can see that for frequencies $0 \leq \omega \leq 2\pi$

rad/sec ($0 \leq f \leq 1$ Hz), the plant transfer function G_p can be approximated with an integrator

$$G_p \cong G_{p0} = (j\omega)^{-1} I_2 \quad \text{for } 0 \leq \omega \leq 2\pi \text{ rad/sec} \quad (4)$$

Thus, the feed-forward differentiator

$$G_f = j\omega I_2 \quad (5)$$

will satisfy Eq. (3) in the frequency range $0 \leq \omega \leq 2\pi$ rad/sec. In Fig. 2(a), the diagonal terms of the differentiator transfer function of Eq. (5) are shown with dotted lines. Their inverses (dashed lines) are equal to the plant transfer function, as in Fig. 2 for frequencies up to 1 Hz. The off-diagonal terms of Eq. (5) (transfer functions from elevation command to azimuth position, and from azimuth command to elevation position) should be zero; actually, they are small for frequencies up to 1 Hz, as in Fig. 2(b) and Fig. 2(d).

The closed-loop transfer functions for a system with and without the feed-forward gain are compared in Fig. 3. Figures 3(a) and 3(c) show that for frequencies up to 1 Hz, the system with the feed-forward gain has superior tracking properties when compared with the system without feed-forward gain. This is confirmed by tracking simulations with a trajectory like that in Figs. 4(a) and

4(b). The DSS-13 antenna, with the proportional gain $k_p = 0.5$, and the integral gain $k_i = 1.8$ in azimuth and elevation, was investigated. The tracking errors in elevation and cross-elevation are compared for the antenna with the feed-forward loop (Fig. 5) and without the feed-forward loop (Fig. 6). A significant improvement in tracking accuracy for the system with the feed-forward loop was observed, namely, from 73.1 to 1.4 mdeg in elevation, and from 60.1 to 0.2 mdeg in cross-elevation. However, the high-frequency components of the command are strongly amplified for the system with feed-forward gain when compared with the system without feed-forward gain. This effect can be observed from the transfer function plots in Fig. 3, where the resonance peaks of the system with feed-forward gain are much higher than the ones of the system without feed-forward gain. Also the intensive oscillatory motion in the pointing error plots (see Fig. 5) is observed. As a result, any sharp change in the command may cause excessive vibrations of the antenna.

Despite the increased sensitivity to the command inputs, the disturbance rejection of the antenna with feed-forward gain remains the same as that for the antenna without feed-forward gain. This follows from Eq. (2), where it is shown that the tracking error e due to wind disturbance w is independent of the feed-forward gain G_f . Thus the pointing errors due to wind gust disturbances are comparable with the results obtained for the DSS-13 antenna with the PI controller (see [2]).

References

- [1] W. Gawronski and J. A. Mellstrom, "Modeling and Simulations of the DSS 13 Antenna Control System," *TDA Progress Report 42-106*, vol. April-June, Jet Propulsion Laboratory, Pasadena, California, pp. 204-248, August 15, 1991.
- [2] W. Gawronski, B. Bienkiewicz, and R. E. Hill, "Pointing-Error Simulations of the DSS-13 Antenna Due to Wind Disturbances," *TDA Progress Report 42-108*, vol. October-December, Jet Propulsion Laboratory, Pasadena, California, pp. 109-134, February 15, 1992.

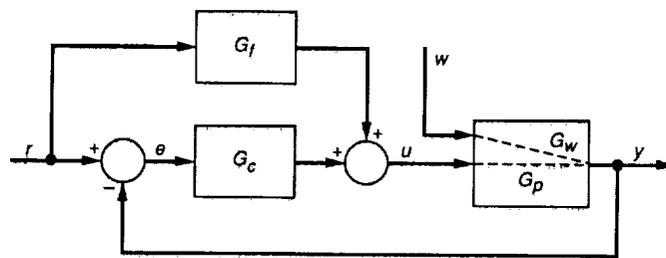


Fig. 1. Antenna control system with the feed-forward loop.

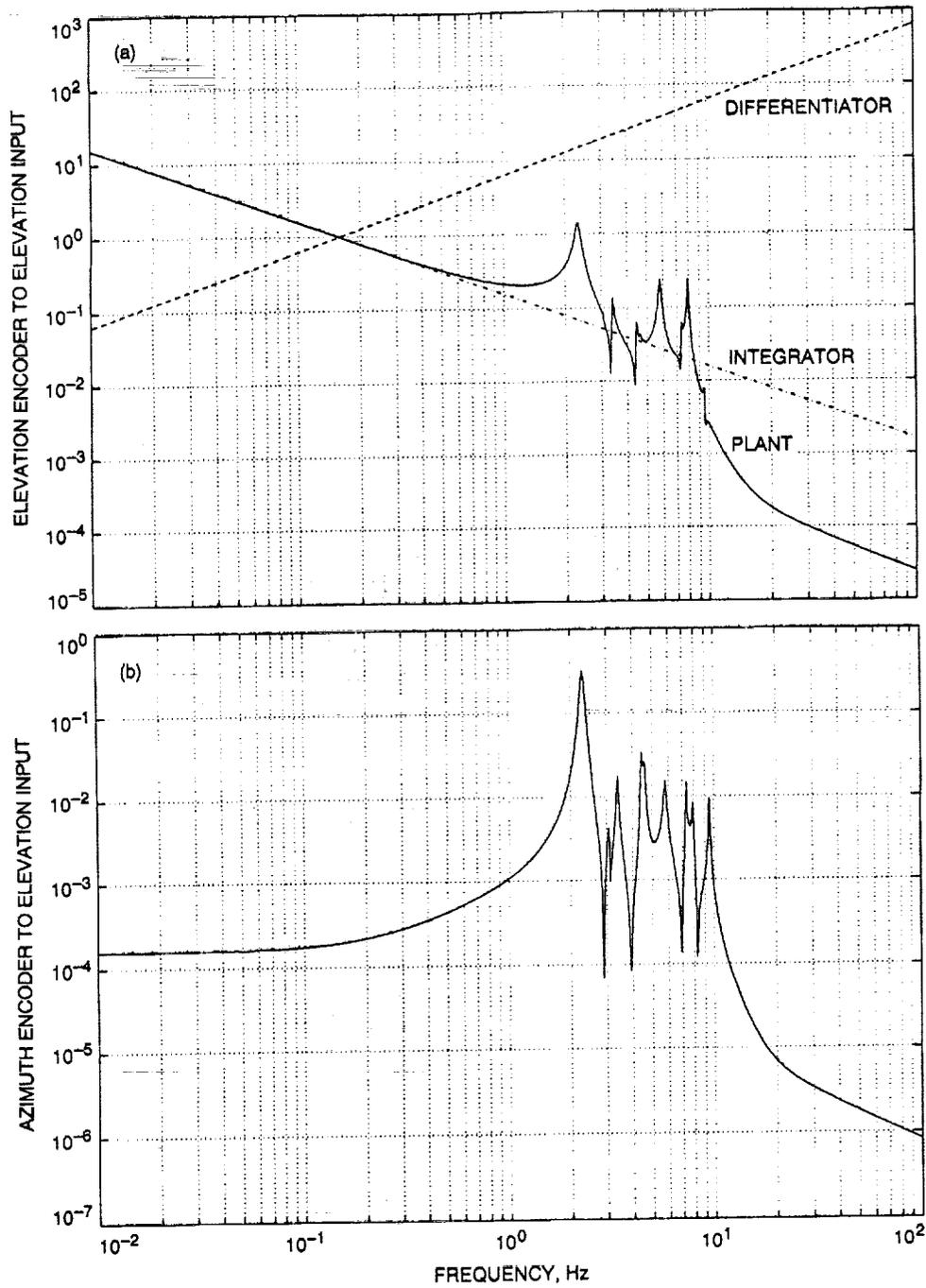


Fig. 2. Transfer functions of antenna rate loop model and of differentiator and integrator: (a) elevation encoder to elevation input; (b) azimuth encoder to elevation input; (c) azimuth encoder to azimuth input; and (d) elevation encoder to azimuth input.

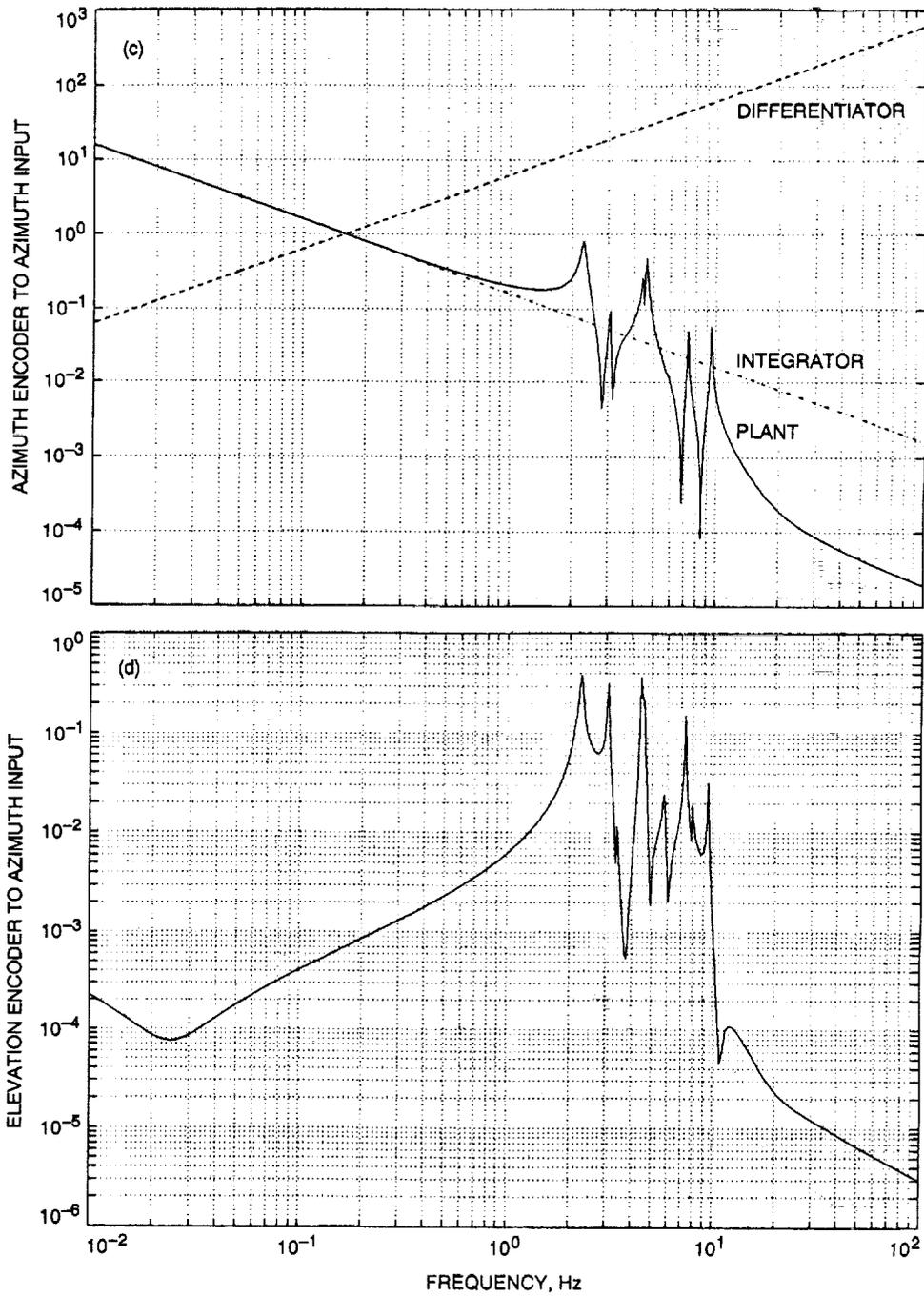


Fig. 2 (contd).

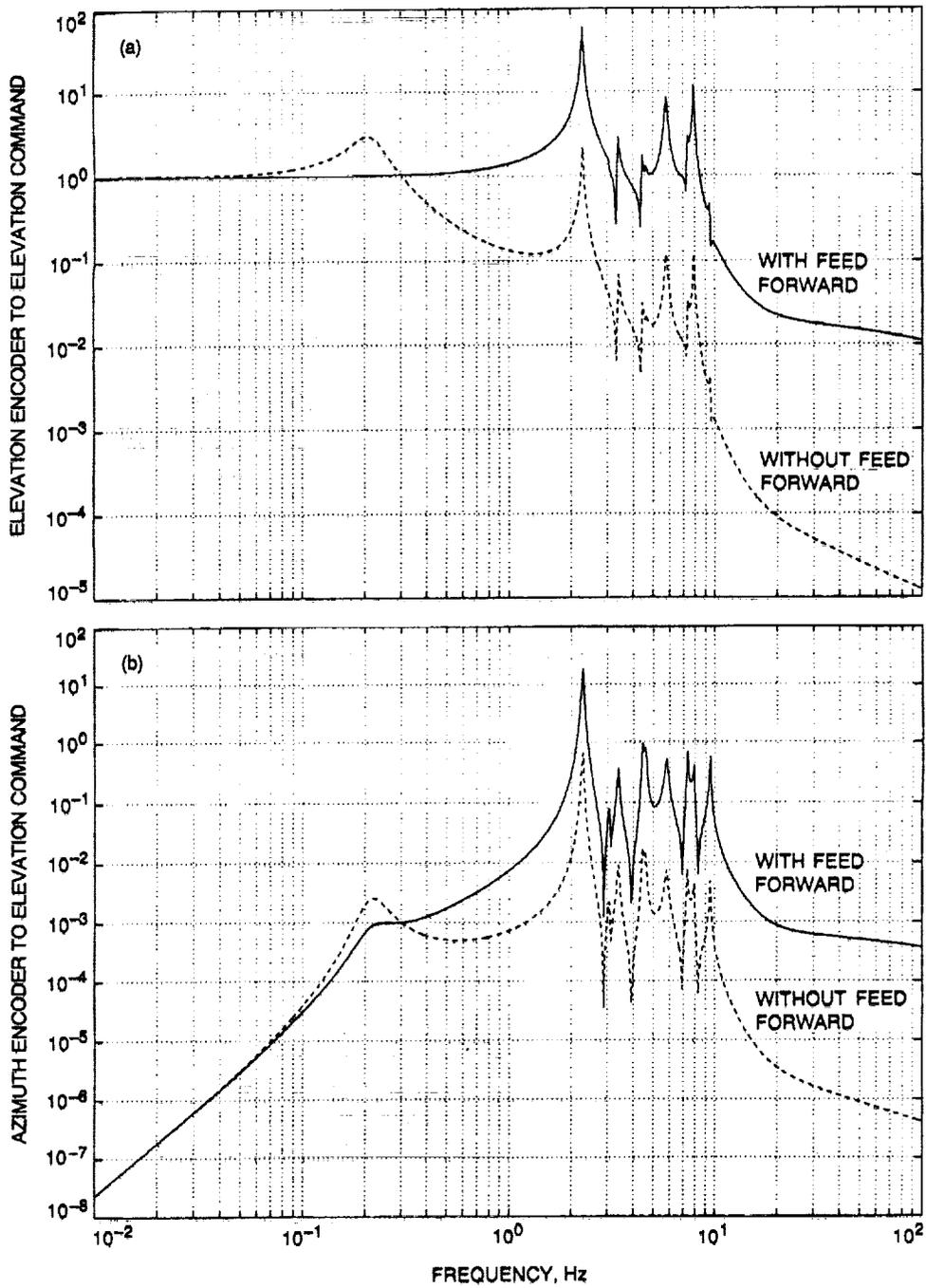


Fig. 3. Closed-loop transfer functions—with feed-forward loop and without feed-forward loop: (a) elevation encoder to elevation command; (b) azimuth encoder to elevation command; (c) azimuth encoder to azimuth command; and (d) elevation encoder to azimuth command.

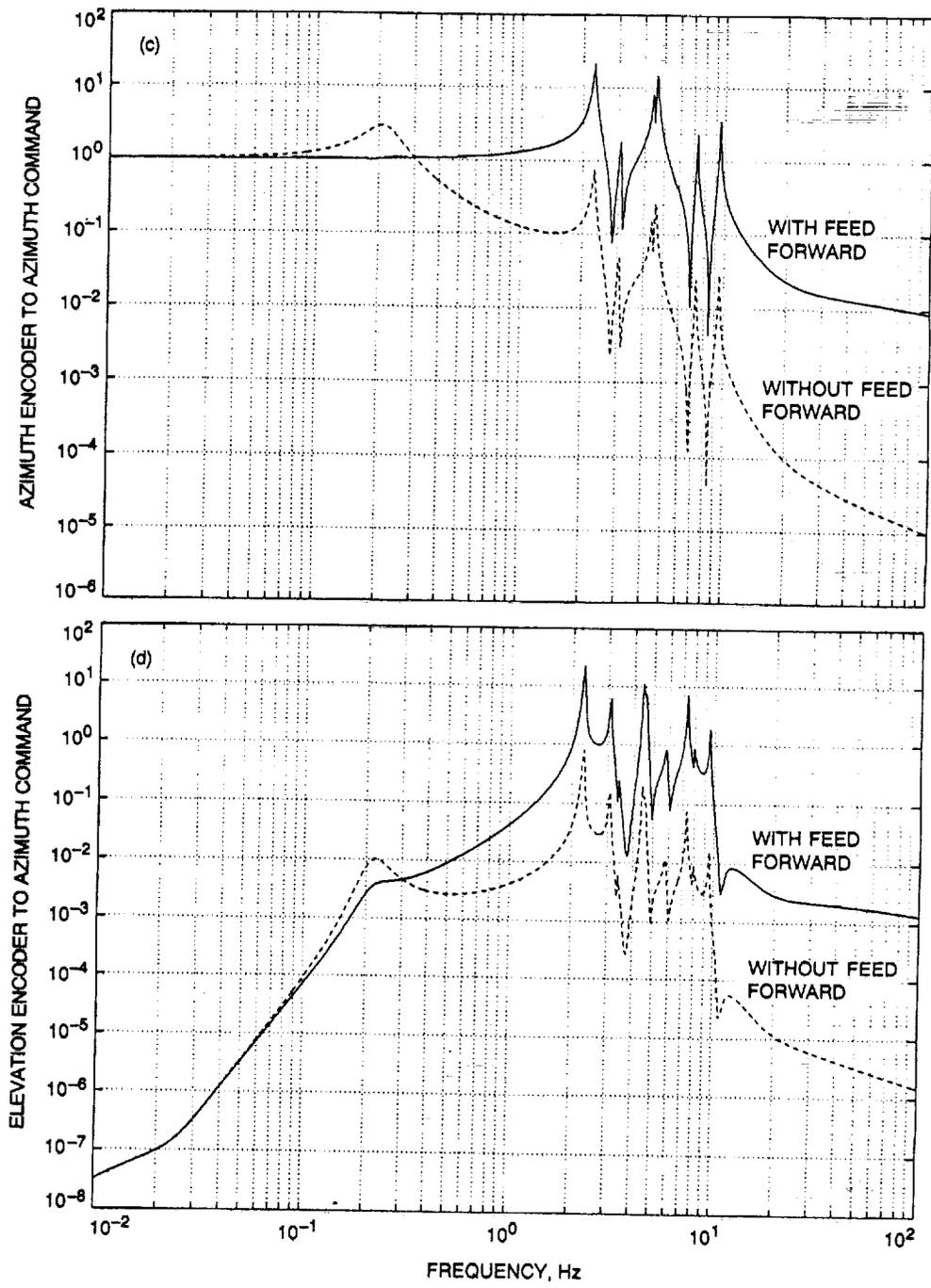


Fig. 3 (contd).

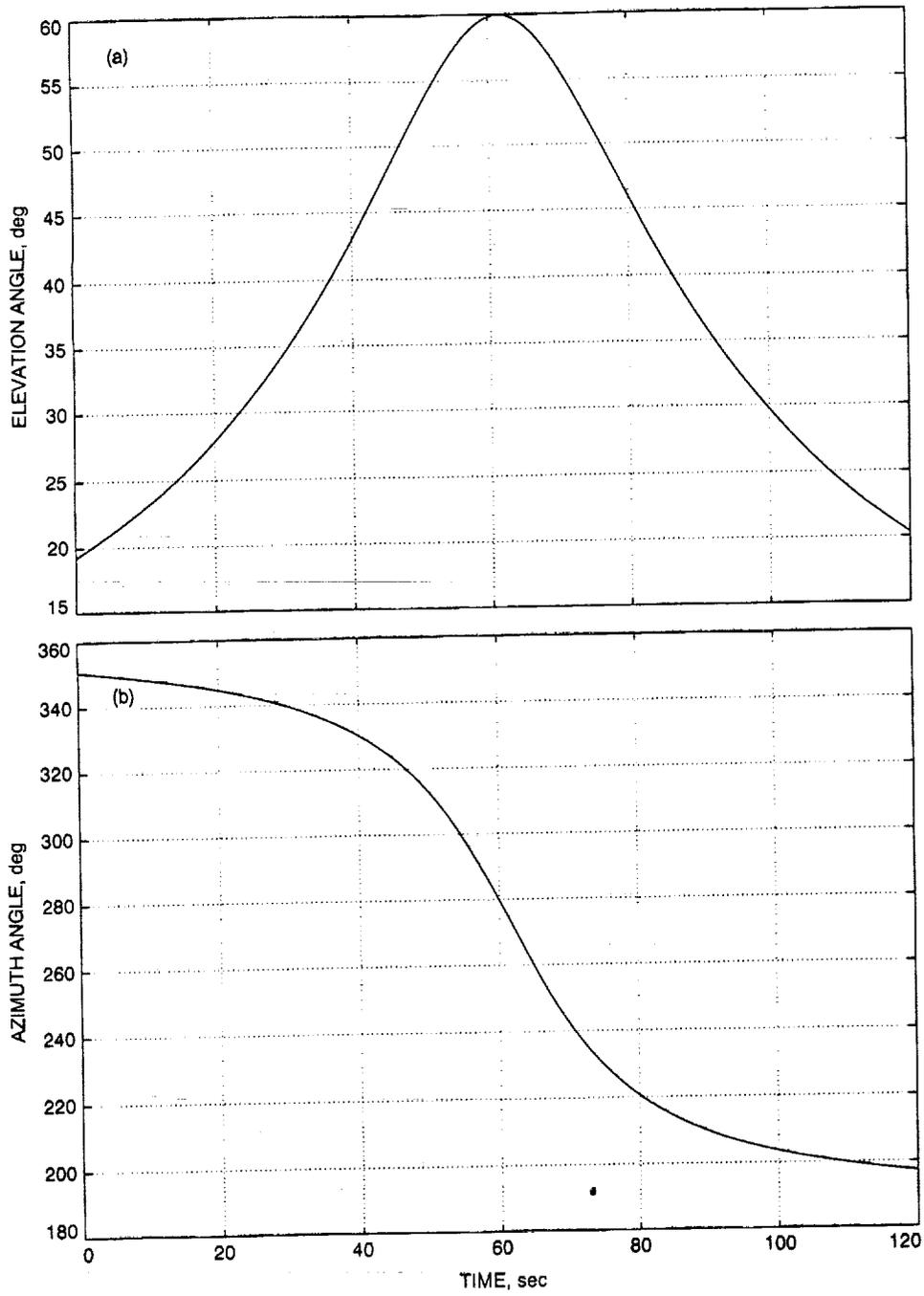


Fig. 4. Trajectory used for simulations: (a) in elevation and (b) azimuth.

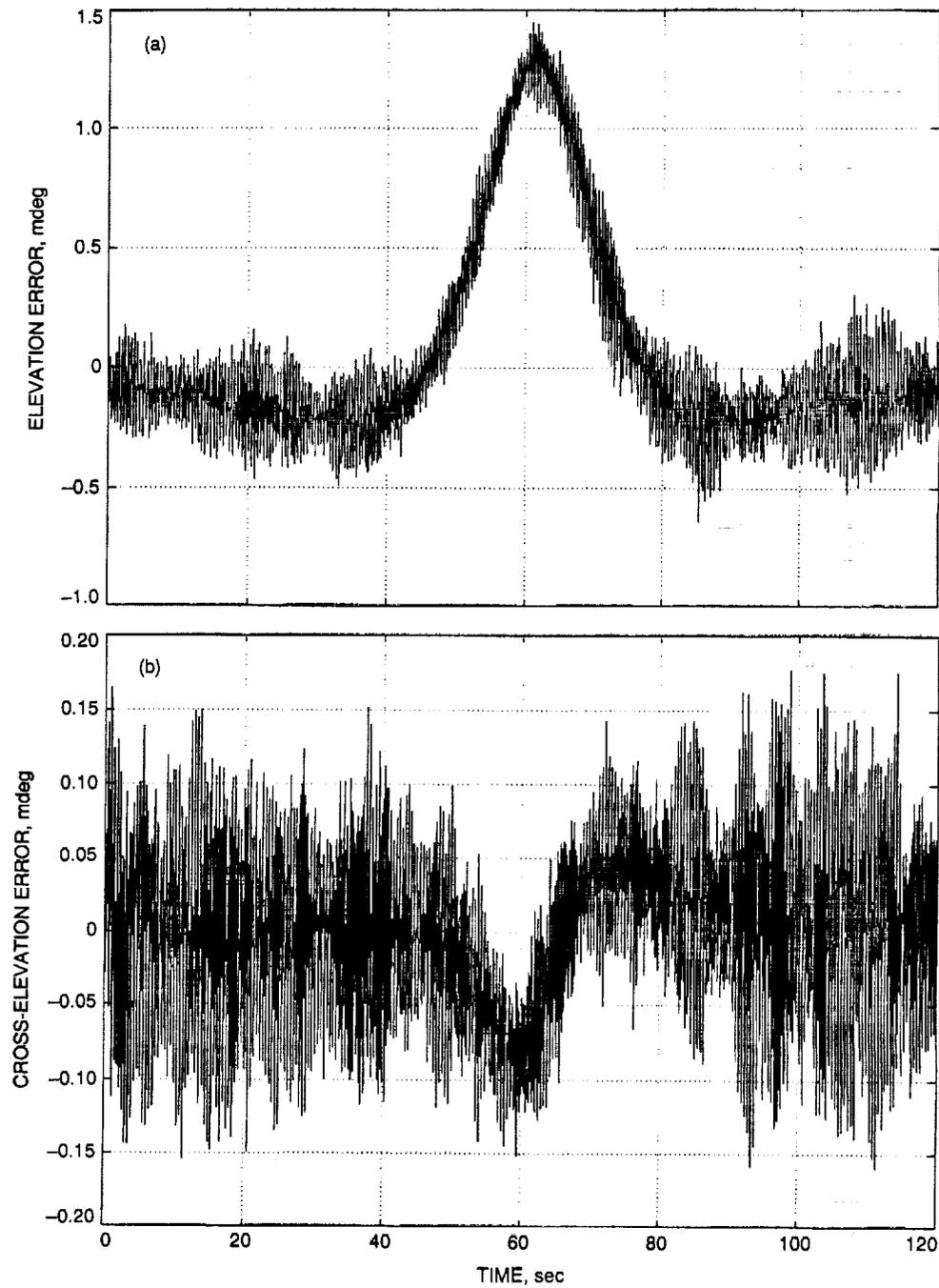


Fig. 5. Pointing errors for the control system with the feed-forward loop: (a) elevation error and (b) cross-elevation error.

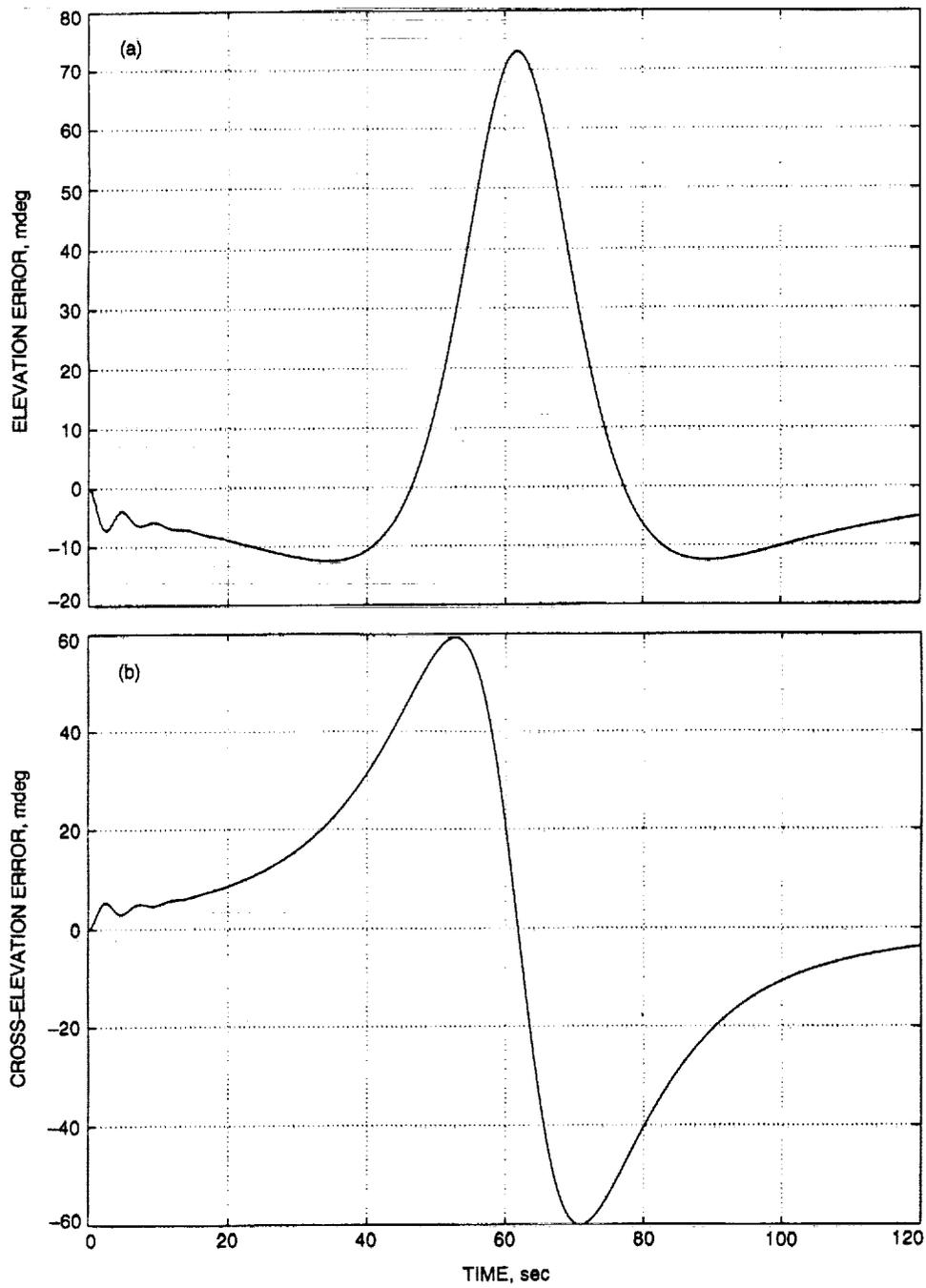


Fig. 6. Pointing errors for the control system without the feed-forward loop: (a) elevation error and (b) cross-elevation error.

N 9 3 - 1 9 4 8 5

120455

9 15

Availability Analysis of the Traveling-Wave Maser Amplifiers in the Deep Space Network

Part I: The 70-Meter Antennas

T. N. Issa

TDA Mission Support and DSN Operations

This article describes the results of the reliability and availability analyses of the individual S- and X-band traveling-wave maser (TWM) assemblies and their operational configurations in the 70-meter antennas of NASA's DSN. For the period 1990 through 1991, the TWM availability parameters for the Telemetry Data System are: mean time between failures (MTBF), 930 hr; mean time to restore service (MTTRS), 1.4 hr; and the average availability, 99.85 percent. In previously published articles, the performance analysis of the TWM assemblies was confined to the determination of the parameters specified above. However, as the mean down time (MDT) for the repair of TWMs increases, the levels of the TWM operational availabilities and MTTRS are adversely affected. In this article, a more comprehensive TWM availability analysis is presented to permit evaluation of both MTBF and MDT effects. Performance analysis of the TWM assemblies, based on their station monthly failure reports, indicates that the TWMs required MTBF and MDT levels of 3000 hr and 36 to 48 hr, respectively, have been achieved by the TWMs only at the Canberra Deep Space Station (DSS 43). The Markov Process technique is employed to develop suitable availability measures for the S- and X-band TWM configurations when each is operated in a two-assembly standby mode. The derived stochastic expressions allow for the evaluation of those configurations' simultaneous availability for the Antenna Microwave Subsystem. The application of these expressions to demonstrate the impact of various levels of TWM maintainability (or MDT) on their configurations' operational availabilities is presented for each of the 70-m antenna stations.

I. Introduction

The traveling-wave maser (TWM) assemblies exist as part of the Antenna Microwave Subsystem (UWV) at each Deep Space Station (DSS). They are used on the large antenna structures at the Goldstone, Madrid, and Canberra communication complexes. The technical performance of TWMs as well as their operational availability are critical

factors in the performance of the DSN. The important performance and operational characteristics of these TWMs are presented in Section II.C of this article.

Previous analyses of the performance reliability and availability of the DSN TWMs are described in [1,2]. These were confined to the mean time between failure

C-4

(MTBF) determination for the TWMs network-wide, and their average availability for the DSN Telemetry System for the period 1981 through 1983. The results indicated that the averaged MTBF (1200 hr) was considerably below the desired level and the averaged MTTRS (2-5 hr) was very long as compared with the requirements. Therefore, a number of recommendations were made to improve the TWM assembly availability characteristics.

Achieving high levels of reliability and availability for the TWM assemblies is a very demanding task. It requires continued improvements in the TWM assembly maintenance procedures, as well as in the assembly's engineering development process. In recent years, it has become apparent that the availability of the TWM assemblies has been degrading. Their operational MTBF is currently believed to be in the range of 1000 to 1400 hr (except for the Canberra complex). This situation demonstrates the need to quantitatively assess the operational availability of the TWMs.

This article is the first of two on the availability modeling and analysis of the DSN TWM configurations. The analysis in this article, Part I, was performed on the TWM configurations in the 70-m antenna subnet. Part II will present the results of the analysis performed on the TWM configurations in the 34-m antenna subnets.

Section II of this article provides the functional description and operational characteristics of these TWMs. Section III presents the application of the Fault Tree Analysis (FTA) technique to the failure analysis of a TWM assembly and provides the results of the quantitative availability evaluation of individual TWMs. Section IV applies a Markov Process technique to derive the steady-state availability expressions for the TWM configurations. Section V describes the application of the derived expressions to the availability analysis of the TWM configurations in the 70-m antennas. Section VI provides a discussion of the study findings and related observations. The last section includes a summary and concluding remarks.

II. Description of TWM Configurations in the 70-M Antennas

A. Functional Configuration

The TWM assemblies are the dominant elements of the UWV at each of the 70-m antenna stations (DSS 14, DSS 43, and DSS 63). There are two TWM assemblies providing low-noise amplification at each of the S-band and X-band frequency feeds. Figure 1 represents a functional

block diagram of these TWM configurations as part of a 70-m Antenna Microwave Subsystem.

A single TWM assembly consists basically of a maser (amplifier) and its closed-cycle refrigerator (CCR), referred to simply as the TWM assembly. Figure 2 shows the TWM and CCR equipment connections.

B. Theory of Operation

The TWM assembly provides low-noise amplification for the received S- or X-band frequency. To accomplish this, the maser (amplifier) is cooled to approximately 4.5 K using closed-cycle refrigeration. To initiate a maser cooldown, the helium compressor (Fig. 2) circulates helium gas to and from the helium refrigerator at controlled pressure and flow rates. The helium refrigerator cools the helium until some of the gas reaches a liquid state (4.5 K) at the bottom of the refrigerator (the cold station). The maser is physically attached to the cold station and is cooled to the correct operating temperature of 4.5 K by metallic conduction.¹

C. Functional Requirements

This section contains data related to the performance requirements and operational characteristics (reliability and availability) of the TWM assemblies in the 70-m antenna subnet.

1. Performance Requirements. The TWM assemblies perform the function of low-noise amplification for the received S- and X-band frequencies. The basic performance requirements for the TWMs at each frequency feed are given in Table 1.

The total performance of the entire antenna system (i.e., Antenna Mechanical and Antenna Microwave Subsystems) is impacted by gain and noise temperature contributions of the TWMs, as well as by other critical elements of the system.

2. Operational Characteristics. The TWM assemblies were designed to demonstrate an operational (field) performance that meets the following characteristics:²

- (1) Functional availability of at least 0.998, since the desired availability for the Antenna Microwave Subsystem is at least 0.996.

¹ General Procedures, "S-Band Traveling-Wave Maser Group Block IV," TM03703 (internal document), Jet Propulsion Laboratory, Pasadena, California, September 15, 1979.

² Communication Complex System Functional Requirements (1991-1996), "Antenna Microwave Subsystem," 824-16 (internal document), Jet Propulsion Laboratory, Pasadena, California, November 1, 1991.

- (2) An MTBF level in excess of 3000 hr.
- (3) MTTRS for a single-channel operation (TWM backup is available) of 5 min. with a maximum of 10 min.
- (4) Mean down time (MDT) of a TWM assembly for repair maintenance is a maximum of 48 hr.

A TWM failure (due to its CCR failure) is recognized by the operator when deviations from requirements are detected in the CCR operating parameters (temperature levels, pressure levels, and helium flow rates). Also, major discrepancies in the observed antenna gain to system noise temperature ratio would indicate a possible TWM failure.

III. Individual Performance Analysis of TWM Assemblies

In this section, the individual performance analysis of TWM assemblies is presented using both qualitative and quantitative approaches. The two basic objectives of this analysis are

- (1) To develop an understanding of the TWM failure characteristics pertinent to various failure modes and their effects, which could be a useful tool for TWM-CCR design reviews (Section III.A).
- (2) To identify relevant operational performance measures related to reliability, maintainability, and availability of the TWM assemblies (Section III.B).

A. Qualitative Failure Analysis

1. Failure Cause-Effect Analysis. The failure analysis of a TWM assembly is primarily concerned with identification of its failure effects, referred to as top events, and determining how these can be caused by individual or combined lower level failures or faults. Standard symbols are then used for developing an FTA to describe failure events and their logical connections in relation to predefined top events (effects). The FTA technique and its applications to reliability assessment are discussed in greater detail in [3,4].

To demonstrate the application of the FTA technique to the TWM-CCR failure analysis, data related to failure modes and fault events of TWMs were collected and organized into some general cause-effect relationships. Using only the DSN Discrepancy Reports (DRs) and the stations' monthly failure reports, only top failure effects and contributing failure events were possible to identify. The logical connections between failure events were difficult to determine based on the current station reporting forms.

Therefore, only a simple FTA is shown in Fig. 3 for a TWM assembly, in which the top event is an operational failure. The FTA shows that this failure may be caused by any of the first-level failure modes (four-input OR gate). The analysis then proceeds, as shown, by determining how each of the lower level failure events can be caused by individual basic faults or events.

It is noteworthy that a different FTA will have to be constructed for each top failure event of the TWM assembly that can be caused by various relations between lower failure events. Therefore, individual FTA should be performed for the major failure effects of the TWM, such as cooldown failure, refrigerator low flow, and compressor failure. These analyses are particularly relevant to the development, operation, and maintenance improvements of the TWM assemblies, since critical levels of failure would indicate their contributions to the TWM reliability and availability.

2. Failure-Cause Distribution. Analysis of 1990-1991 TWM station failure reports gave the probability distribution of the major areas (causes) contributing to the 70-m TWM assembly failure. The distribution is summarized in Table 2. Refrigerator contamination and low helium flow proportions constitute about 65 percent of total TWM failures. More than 75 percent of the refrigerator failures were attributed to contamination of the helium gas. A significant proportion of TWM failure was attributed to the helium compressor failure.

B. Quantitative Availability Analysis

In this section, relevant availability factors and measures are quantitatively determined for the individual TWM assemblies in the 70-m antenna subnet. The factors include MTBF, MDT for corrective maintenance, and mean preventive-maintenance time (MPT); the measures include operational availability (A_o) and achieved availability (A_a).

Data related to TWM assembly outages and service restorations (during spacecraft scheduled support periods) were taken from the DR System. However, since TWM assembly failures (basically hardware-fault related) occurring over nonspacecraft support periods are not reported by the DR system, it is thought that the station monthly failure and maintenance reports should also be reviewed for this analysis.

1. Reliability (MTBF) Analysis. Earlier analysis of the TWM failure data based on the DSN DRs was reported in [1,2]. Network-wide estimates of MTBF as reported for three different periods are shown in Table 3. An

operational MTBF level of approximately 2500 to 3000 hr was considered a desirable and achievable reliability target for individual TWM assemblies.

This reliability analysis considers the individual performances of the TWMs at DSS 14, DSS 43, and DSS 63. The MTBF analysis is based on the station monthly failure reports for the period from November 1990 through December 1991. Estimated MTBF levels for the individual TWMs at each station are summarized in Table 3. An MTBF level of 2400 to 3000 hr has been achieved by the TWMs at the Canberra Complex (DSS 43). The TWMs at DSS 14 demonstrated an MTBF level in the range of 1400 to 1600 hr, whereas the MTBF level for the TWMs at DSS 63 is in the range of 1400 to 1900 hr.

2. Availability Analysis. Determination of TWM availability includes (1) a data-availability estimate of TWMs relative to each DSN data system, which is determined for a total scheduled mission-support time over a given period and (2) an overall-availability estimate of TWMs determined for the total period considered (including mission and other activity support times).

a. Data-Availability of TWMs. This availability for any data system is a function of both MTBF (or data-outage rate) and MTTRS. Separate analyses were conducted on the Antenna Microwave Subsystem availability for the Telemetry Data System and the individual 70-m station TWM availability for the Telemetry Data System for the period 1990 through 1991. The outage data for these analyses were obtained from the DR system, and the results are shown in Tables 4 and 5, respectively.

The average TWM availability for telemetry data is 99.8 percent. The average MTBF of a TWM assembly (approximately 930 hr) is good relative to that of other major Telemetry System elements; however, it is considerably lower than the desired MTBF level of 2500 to 3000 hr. The average MTTRS (1.4 hr) is relatively high and far off the service restoration requirements (these include mean and maximum durations of 15 and 30 min, respectively, for all support activities).

b. Overall availability of TWMs. This availability includes the determination of the TWM operational availability (A_o) and achieved availability (A_a) levels. For evaluating these measures, failure and maintenance data were taken from the 70-m stations' TWM monthly performance reports. Then operational and achieved availabilities for the individual TWMs were determined, as shown in Table 6, for the period November 1, 1990, through December 31, 1991, using the relationships provided in [3,4] as follows:

$$A_o = \text{MTBF}/(\text{MTBF}+\text{MDT}) \quad (1)$$

$$A_a = \text{MTBF}/(\text{MTBF}+\text{MDT}+\text{MPT}) \quad (2)$$

where

MTBF = MTTF (mean time to failure only)

MDT = (waiting-time to restore service
+ logistic delay time
+ corrective maintenance time
+ CCR decontamination time
+ assembly cooldown and testing time)

MPT = (preventive-maintenance time
+ decontamination time
+ cooldown and testing time)

The lower MDT levels for DSS-43 TWM assemblies indicate an improvement in both repair and operating procedures. The TWMs at both DSS 14 and DSS 63 demonstrated relatively higher MDT levels.

IV. Operational Availability Modeling of TWM Configurations

In this section, the operational availability performance of the X- and S-band TWM configurations is being considered for modeling. Each configuration is most frequently operated as a two-assembly standby configuration. The availability measure for the case of the two-assembly parallel configuration was derived in [5,6].

The purpose of this modeling effort is to derive the steady-state availability measures for the TWM standby configuration using the Markov Process technique discussed in [6,7]. In the following subsections, the model assumptions, model formulation, and the derived availability measures for the TWM standby configuration are presented.

A. Model Assumptions

The following are the basic assumptions of the Markov model for the TWM configurations under study:

- (1) The model represents an S- or X-band TWM standby configuration consisting of two redundant assemblies and a single maintenance technician.

- (2) The stochastic failure and repair processes for the operating assembly are stationary (constant failure and repair rates).
- (3) The repair activity and its duration (total downtime) of the failed TWM covers the time for corrective maintenance and the total time for decontamination and cooldown processes.
- (4) The backup TWM is warm (power-connected) but in a nonoperating mode (i.e., it receives only a duplicate signal and does not provide output to other processors in the link). Therefore, it is assumed that the standby TWM has a negligibly small failure rate (an approximate zero rate) relative to that of the prime operating TWM.
- (5) In the initial operating state of the TWM configuration, State 1, one of the TWM assemblies is in an operable standby mode.
- (6) During the TWM operation, only one change can take place in the state of the configuration at each instantaneous increment of time.

B. Model Formulation

The state-space diagram associated with the TWM standby configuration under study is shown in Fig. 4. The following symbols are associated with this diagram:

i denotes the i th state of the TWM configuration, for $i = 1, 2, 3$; where $i = 1$ (one assembly is operating, the other assembly is in standby mode); $i = 2$ (one assembly failed and is under repair, the other assembly is linked and operating); $i = 3$ (both TWMs are down, and one of them is under repair).

P_i denotes the steady-state probability that the TWM configuration is in State i , for $i = 1, 2, 3$.

λ denotes the constant failure rate of the operating TWM assembly.

μ denotes the constant repair and maintenance of the failed TWM assembly.

The steady-state availability model for this configuration is developed using the frequency-balance principle of the Markov Process theory as described in [5,7]. The frequency balance equations for the three-state availability model of Fig. 4 can be written as

$$\text{State 1: } \lambda P_1 = \mu P_2 \quad (3)$$

$$\text{State 2: } (\lambda + \mu) P_2 = \lambda P_1 + \mu P_3 \quad (4)$$

$$\text{State 3: } \mu P_3 = \lambda P_2 \quad (5)$$

Using Eqs. (3) and (5) and the unity equation $P_1 + P_2 + P_3 = 1$, the state probability expressions are defined as follows:

$$P_1 = \mu^2 / (\lambda^2 + \lambda\mu + \mu^2) \quad (6)$$

$$P_2 = \lambda\mu / (\lambda^2 + \lambda\mu + \mu^2) \quad (7)$$

$$P_3 = \lambda^2 / (\lambda^2 + \lambda\mu + \mu^2) \quad (8)$$

C. The Availability Measures

In the state-space model shown in Fig. 4, States 1 and 2 represent the operating (up) states and State 3 represents the failed (down) state of the configuration. Thus, using the state probabilities given in Eqs. (6)–(8), the steady-state operational availability measure of the standby configuration, denoted by A_o , is given by

$$A_o = (\mu^2 + \lambda\mu) / (\lambda^2 + \lambda\mu + \mu^2) \quad (9)$$

and the steady-state unavailability measure of the standby configuration, denoted by U_o , is given by

$$U_o = \lambda^2 / (\lambda^2 + \lambda\mu + \mu^2) \quad (10)$$

The availability measure given in Eq. (9) can also be used to develop a simultaneous availability measure of X- and S-band configurations for the UUV operational support at each 70-m antenna station. This is defined as

$$A_o(\text{for antenna microwave}) = A_o(\text{X-band}) \times A_o(\text{S-band}) \quad (11)$$

The application of Eqs. (9) and (11) to evaluate the availability of the X- and S-band TWM configurations at each 70-m antenna station is presented in Section V.

V. Applications and Analysis

For the application of the availability measures given in Eqs. (9) and (11), pooled MTBF and MDT levels were estimated for the X- and S-band TWM configurations using their individual assembly MTBF and MDT levels, which were computed earlier and listed in Table 6. The pooled MTBF estimate for a TWM was computed as the average of MTBF levels of the individual assemblies (in a configuration) reduced by 30 to 35 percent of the total variation

of an assembly MTBF level from that average MTBF. On the other hand, the pooled MDT estimate for a TWM was computed as the average of MDT levels of the individual assemblies (in a configuration) rounded off to the closest integer representing an average number of 12-hr maintenance cycles (shifts) required for a TWM repair completion.

The estimates of pooled MTBF and MDT for the TWM assemblies and their use for availability evaluation of the X- and S-band configurations at each 70-m antenna station are described as follows.

A. Availability of TWMs at DSS 14

The estimates of pooled MTBF and MDT for the prime and backup assemblies of the X-band TWM configuration are

$$\text{MTBF} = 1615 \text{ hr, or } \lambda = 0.000619 \text{ failure/hr}$$

$$\text{MDT} = 84 \text{ hr, or } \mu = 0.011904 \text{ repair/hr}$$

and the estimates for each TWM assembly in the S-band configuration are

$$\text{MTBF} = 1560 \text{ hr, or } \lambda = 0.000641 \text{ failure/hr}$$

$$\text{MDT} = 84 \text{ hr, or } \mu = 0.011904 \text{ repair/hr}$$

The application of the availability measure given in Eq. (9) to the TWM configurations at this station results in

$$A_o(\text{X-band configuration}) = 0.997435$$

$$A_o(\text{S-band configuration}) = 0.997256$$

The application of the availability measure given in Eq. (11) results in the following simultaneous availability of TWM configurations for the UWV subsystem:

$$\begin{aligned} A_o(\text{for antenna microwave}) &= A_o(\text{X-band}) \times A_o(\text{S-band}) \\ &= 0.9947 \end{aligned}$$

The simultaneous availability of TWM configurations, which is 0.9947, is slightly lower than its predicted requirement (i.e., 0.996) for the antenna microwave. The

operational availability of each individual configuration is approximately 0.997, as compared with a desirable level of 0.998. This deviation is primarily attributed to both a low MTBF level (1500 to 1600 hr) and a considerable MDT of 84 hr (3 to 3.5 days), as compared with the corresponding parameter levels for the TWMs at the other stations.

The impact of improved TWM assembly reliability on a configuration's operational availability at various levels of assembly MDT is demonstrated in Fig. 5(a). An examination of the plots in this figure indicates that the operational availability requirement for each TWM configuration, which is 0.998, can be achieved or even exceeded when the reliability (MTBF) and maintainability (MDT) parameters of its assemblies meet any of the following practically feasible combinations:

- (1) $\text{MTBF} \geq 1500 \text{ hr; MDT} \leq 60 \text{ hr}$
- (2) $\text{MTBF} \geq 2000 \text{ hr; MDT} \leq 84 \text{ hr}$

B. Availability of TWMs at DSS 43

The estimates of pooled MTBF and MDT for each TWM assembly in the X-band configuration are

$$\text{MTBF} = 2350 \text{ hr, or } \lambda = 0.000425 \text{ failure/hr}$$

$$\text{MDT} = 36 \text{ hr, or } \mu = 0.027777 \text{ repair/hr}$$

and the estimates for each TWM assembly in the S-band configuration are

$$\text{MTBF} = 5080 \text{ hr, or } \lambda = 0.000196 \text{ failure/hr}$$

$$\text{MDT} = 36 \text{ hr, or } \mu = 0.027777 \text{ repair/hr}$$

The application of the availability measures given in Eqs. (9) and (11) to the TWM configurations at this station results in the following individual and simultaneous availabilities:

$$A_o(\text{X-band configuration}) = 0.999768$$

$$A_o(\text{S-band configuration}) = 0.999950$$

$$A_o(\text{for antenna microwave}) = 0.9997$$

The operational availability of both TWM configurations at this station, which is 0.9997, is greater than the

predicted requirement of 0.996 for the period 1990 through 1991. The availability of individual configurations is relatively high and in the range of 0.9997 to 0.9999 for this period. This is primarily attributed to both a reasonable MTBF level (2300 to 2500 hrs) and an acceptable MDT of 36 hours (or 1.5 to 2 days).

The impact of improved TWM assembly reliability on a configuration's operational availability at various levels of assembly MDT is described in Fig. 5(b). The plots demonstrate that the operational availability for each TWM configuration will always meet or exceed the requirement at the current MTBF level (greater than 2000 hr) for as long as the MDT level remains at or below 60 hr.

C. Availability of TWMs at DSS 63

The estimates of pooled MTBF and MDT for each TWM assembly in the X-band configuration are

$$\text{MTBF} = 1820 \text{ hr, or } \lambda = 0.000549 \text{ failure/hr}$$

$$\text{MDT} = 60 \text{ hr, or } \mu = 0.016666 \text{ repair/hr}$$

and the estimates for each TWM assembly in the S-band configuration are

$$\text{MTBF} = 2525 \text{ hr, or } \lambda = 0.000396 \text{ failure/hr}$$

$$\text{MDT} = 36 \text{ hr, or } \mu = 0.027777 \text{ repair/hr}$$

The application of the availability measures given in Eqs. (9) and (11) to the TWM configurations at this station results in the following individual and simultaneous availabilities:

$$A_o \text{ (X-band configuration)} = 0.998948$$

$$A_o \text{ (S-band configuration)} = 0.999799$$

$$A_o \text{ (for antenna microwave)} = 0.9987$$

The operational availability of both TWM configurations at this station, which is 0.9987, is slightly higher than the predicted requirement of 0.996 for the period 1990 through 1991. The availability of the X-band TWM configuration exceeds the required level of 0.998. For the S-band TWM configuration, the availability is relatively high for this period as a result of both a reasonable MTBF level (2525

hours) and an acceptable MDT level of 36 hours (less than 2 days).

The impact of increased TWM assembly reliability on a configuration's operational availability at possible levels of assembly MDT is described in Fig. 5(c). The plots demonstrate that the operational availability for each TWM configuration will always exceed a level of 0.998 at the current MTBF level (greater than 2000 hrs) for as long as the MDT level remains at or below 48 hours.

VI. Discussion of the Findings and Observations

In this section, the results of the operational reliability and availability analysis of the TWM configurations are discussed. The appropriate items related to the improvement of the TWM assembly availability characteristics are also presented.

A. Discussion of the Findings

- (1) Failure analysis of TWM assemblies shows that the helium refrigerator contamination is the dominant cause of TWM failure (55 to 65 percent of total failures). An earlier investigation of this problem implied the feasibility of detecting the development of contamination 1 or 2 days in advance of a TWM failure.
- (2) The proportions of compressor and refrigerator drive-unit failures (12 and 15 percent, respectively) are considered significant, as compared with the failure percentages for other TWM elements. If their failure frequencies were reduced, the TWM MTBF would be considerably improved.
- (3) The DR data show that the average MTTRS for the TWMs (1.4 hr) relative to the DSN Telemetry Data System is about two to three times as long as the MTTRS of other telemetry elements. Improving the MTTRS of the TWM assemblies to a desirable level of 0.7 hr would essentially require improved TWM backup availability. This is accomplished when repair durations, and consequently MDT levels of the TWMs, are reduced.
- (4) Analysis of the 1990-1991 TWM station failure data shows a significant variation in the MTBF level at different complexes. The higher MTBF levels for the TWMs at DSS 43 appear to have resulted from improved repair procedures as well as an increased level of preventive-maintenance activity prior to early 1990. On the other hand, TWMs at DSS 14 and DSS 63 have demonstrated lower MTBF

levels, which are believed to have resulted from less efficient repair and preventive maintenance processes at these stations.

- (5) Analysis of the 1990-1991 TWM station repair data shows a variation in the level of MDT for repair at different complexes. For the TWM assemblies at DSS 43, the MDT level is relatively good (36 hr) and lower than MDT levels demonstrated at the other stations (84 hr at DSS 14 and 48-60 hrs at DSS 63). That is perhaps indicative of improved repair practices as well as better trained maintenance personnel at the Canberra complex.
- (6) TWM operational availabilities at the individual stations (computed by using the proposed measures), are in close agreement with their average availabilities for the Telemetry Data System (generated from the DR system). These availabilities at the individual stations are compared in Table 7.

The variations in corresponding TWM availabilities at different stations are primarily attributed to the effect of incorporating the MDT parameter into the proposed TWM availability measures. The use of TWM MDT is more appropriate than the MTTRS parameter for their operational availability evaluations since the former is more representative of TWM actual unavailable (repair) times.

B. Observations

The following observations, which are primarily drawn from the previous work reported in [2], are based on the findings of the TWM failure-cause analysis.

- (1) Solution of the helium gas contamination problem would reap the greatest dividends. The composition and sources of contamination have to be better understood. Some of the items related to this area include
 - (a) Implementation of the previously proposed computer-based monitoring and data-collection and analysis system for the TWM CCRs would improve the identification of possible CCR faults before they cause TWM failure.
 - (b) Improved field techniques for measuring gaseous impurities and other contaminants in helium would be very valuable.
 - (c) Improved gas flow meters and gauges of the helium refrigerators and compressors would allow for the detection of helium low flow and improve the quality of recorded performance data.

- (2) An understanding of the contribution of the compressor oil to contamination at increased temperatures is needed, as is a mechanism to replace in-service compressor filters and adsorbers without affecting maser operation.
- (3) A detailed failure mode analysis, preferably using the FTA technique, for the helium compressors and refrigerators to reduce their current significant contributions to TWM failure is yet to be done and would be part of another phase.

The following observations are based primarily on the findings of the reliability and availability analyses of the TWMs:

- (1) The preventive-maintenance schedules and procedures for the TWM equipment at both DSS 14 and DSS 63 need to be reviewed to achieve higher MTBF levels for their TWM assemblies.
- (2) The TWM corrective-maintenance procedures and support equipment at DSS 14 and DSS 63 need to be evaluated for possible improvements in order to achieve reduced repair durations and improved TWM backup availabilities. This is essential to reduce the average MTTRS for the TWMs network-wide.
- (3) Improving and sustaining training for the maser operation and maintenance personnel would help to achieve a TWM availability performance consistent with the specified requirements.

VII. Summary

This article has reported the results of the reliability and availability analyses of the TWM assemblies at the 70-m antennas and has presented a stochastic availability evaluation model of their operational configurations. The dominant cause of TWM failures is contamination of the helium gas in the CCRs. This is consistent with the findings of a previous study; however, another important finding is that proportions of TWM failures attributed to compressor and refrigerator drive-unit failures have almost doubled in recent years. The current MTBF level of the TWMs for spacecraft support (approximately 930 hr) can practically be improved. The average MTTRS can be reduced to 0.7 or 0.8 hr if TWM backup-assembly availability is improved.

The MTBF and MDT levels of the TWMs at the Canberra Complex DSS 43 indicate that it is possible to achieve the desired levels of these parameters (MTBF of

2500-3000 hr; MDT of 36-48 hr) for the present TWMs at both DSS 14 and DSS 63 when operation and maintenance procedures are consistently improved and practiced. Individual availabilities of TWMs at DSS 14 are considerably lower than corresponding levels for the TWMs at DSS 43 and DSS 63. That is indicative of the adverse impact of the relatively higher MDT (84 hr) for the TWMs at DSS 14.

The derived stochastic expressions provide adequate measures of the S- and X-band TWM standby configura-

tion availabilities and allow for the evaluation of simultaneous operational availability of these configurations for their Antenna Microwave Subsystem. The application of these availability expressions to the 70-m antenna TWM configurations indicates a relatively lower operational availability level achieved at DSS 14. That is primarily attributed to the higher MDT of the failed TWM at this station. The proposed measures can be considered as useful tools to examine possible MTBF and MDT trade-offs that would result in an improved TWM configurations' operational availabilities at the 70-m antennas.

Acknowledgments

The author thanks H. Matossian and M. Loria for their assistance in the review and analysis of station failure and maintenance reports; D. Custer for providing summaries of Discrepancy Report data; R. Fuzie and G. Rogers for their useful discussions and review of this article; and D. Trowbridge for his careful review of a draft of this article. The TWM maintenance engineering and support staff at the Goldstone, Canberra, and Madrid Communication Complexes are acknowledged for their efforts in maintaining the maser performance databases as well as their support of this study.

References

- [1] R. Stevens, "Availability of the DSN Telemetry Data System and Its Major Elements, Including the TWM Assemblies," *TDA Progress Report 42-78*, vol. April-June 1984, Jet Propulsion Laboratory, Pasadena, California, pp. 184-191, August 15, 1984.
- [2] R. Stevens and C. P. Wiggins, "A Study of DSN Traveling-Wave Maser System Reliability," *TDA Progress Report 42-78*, vol. April-June 1984, Jet Propulsion Laboratory, Pasadena, California, pp. 192-198, August 15, 1984.
- [3] P. D. T. O'Connor, *Practical Reliability Engineering*, 2nd ed., New York: John Wiley and Sons, 1985.
- [4] E. J. Henley and H. Kumamoto, *Probabilistic Risk Assessment: Reliability Engineering, Design, and Analysis*, New York: The Institute of Electrical and Electronics Engineers, Inc., 1991.
- [5] R. Billinton and R. N. Allan, *Reliability Evaluation of Engineering Systems: Concepts and Techniques*, New York: Plenum Press, 1983.
- [6] G. H. Sandler, *System Reliability Engineering*, Englewood Cliffs, New Jersey: Prentice-Hall, Inc., 1963.
- [7] E. Cinlar, *Introduction to Stochastic Processes*, Englewood Cliffs, New Jersey: Prentice-Hall, Inc., 1975.

Table 1. Basic performance requirements of the 70-m antenna TWM assemblies.

Performance parameter	S-band (Block III, IV, V)	X-band (Block IIA)
Frequency, MHz	2270-2300	8400-8500
Gain contribution, dB		
Peak	45	45
Minimum	44	44
Gain stability, dB		
Stationary-10 sec	0.03-0.05	0.03
Stationary-12 hr	0.5	0.5
Moving	0.5-2.0	0.5
Noise temperature contributions, K	8-10 Block III <5.0 Blocks IV, V	<4.0

Table 2. Failure-cause distribution of the 70-m antenna TWM assemblies.

Area of failure	Number of TWM failures			Number of failures (subnet)	Total failures, percent
	DSS 14	DSS 43	DSS 63		
Helium refrigerator	24	16	25	65	69.8
Contamination	(19)	(12)	(20)	(51)	(54.8)
Other (drive unit, etc.)	(5)	(4)	(5)	(14)	(15.0)
Helium compressor	4	2	5	11	11.8
Low flow	4	2	3	9	9.7
Power supply/distribution	2	1	2	5	5.4
Miscellaneous (pump, klystron, maser, etc.)	0	2	1	3	3.3
Total	34	23	36	93	100

Table 3. MTBF history of the DSN TWM assemblies.

Time period	Approximate MTBF, hr	Remarks
DSN DR's data for		
Late 1960s	1000-1300	All complexes/stations
Late 1970s (1979-1981)	3000	All complexes/stations
Early 1980s (1982-1983)	1000-1200	All complexes/stations
Station reports' data for		
Early 1990s (1990-1991)	1400-1600 (S-band) 1600 (X-band) 5000 (S-band)	Goldstone Complex; only DSS14 Canberra Complex; only DSS 43
	2400-4000 (X-band) 1900-2500 (S-band) 1400-1800 (X-band)	Madrid Complex; only DSS 63

Table 4. Telemetry data system availability based on its major subsystem contributions.

Subsystem	Number of outages	Total outage, hr	MTBF, hr	MTTRS, hr	Telemetry availability, percent ^a
Antenna Mechanical	229	200	171	0.9	99.50
Antenna Microwave (TWM included)	69	96	567	1.4	99.76
Radio Frequency Interference	00	160	244	0.9	99.61
Receiver	195	74	220	0.4	99.82
Telemetry	472	271	92	0.6	99.35
Facility	34	35	3760	0.6	99.94

^aData availability = MTBF/(MTBF + MTTRS)

Table 5. Telemetry data availability analysis based on the contributions of the 70-m station TWMs.

Station TWMs	Schedule support time, hr	Number of outages	Total outage, hr	MTBF, hr	MTTRS, hr	Telemetry availability, percent ^a
DSS 14	8356	15	14.6	557	0.97	99.83
DSS 43	9165	1	8.3	9165	8.3	99.91
DSS 63	9533	13	17.4	734	1.3	99.82
Total/average	27054	29	40.3	933	1.4	99.85

^aData availability = MTBF/(MTBF + MTTRS)

Table 6. Availability characteristics of TWM assemblies in the 70-m antenna subnet.

DSS complex	TWM type	Availability factor			Availability measure	
		MTBF ^a	MDT ^b	MPT ^c	A _o	A _a
DSS 14, Goldstone, California	S1	1361	92.6	48	0.93629	0.90636
	S2	1963	72.0	48	0.96462	0.94239
	X1	1628	68.0	48	0.95990	0.93348
	X2	1601	95.0	48	0.94399	0.91800
DSS 43, Canberra, Australia	S1	5070	42.0	N/A ^d	0.99178	0.99178
	S2	5091	21.0	N/A	0.99589	0.99589
	X1	5079	33.0	N/A	0.99354	0.99354
	X2	1232	46.75	N/A	0.96344	0.96344
DSS 63, Madrid, Spain	S1	1987	38.4	32	0.98104	0.96578
	S2	5056	20.0	36	0.99605	0.98905
	X1	1070	54.0	30	0.95196	0.92721
	X2	3568	64.0	36	0.98238	0.97274

^a The MTBF levels shown are calculated as the total operating hours of each TWM divided by the number of TWM failure (outage) events for the period considered.

^b The MDT shown is calculated as the total outage times and repair times divided by the number of reported failures for each TWM assembly.

^c The MPT shown is calculated as the total preventive-maintenance time divided by the number of preventive activities performed for a period.

^d Indicates that preventive-maintenance activities were not scheduled for the period considered.

Table 7. A comparison of TWM availabilities at the 70-m stations.

DSS-TWMs	Operational availability (proposed measures), percent	Average availability (DSN DR system), percent
DSS 14	99.73-99.75	99.83
DSS 43	99.97-99.99	99.91
DSS 63	99.89-99.97	99.82

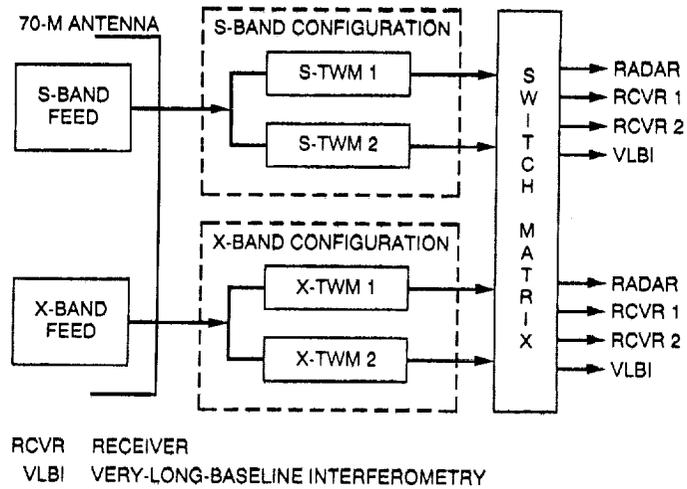


Fig. 1. Functional block diagram for the traveling wave maser (TWM) amplifiers in the 70-m antenna microwave subsystem.

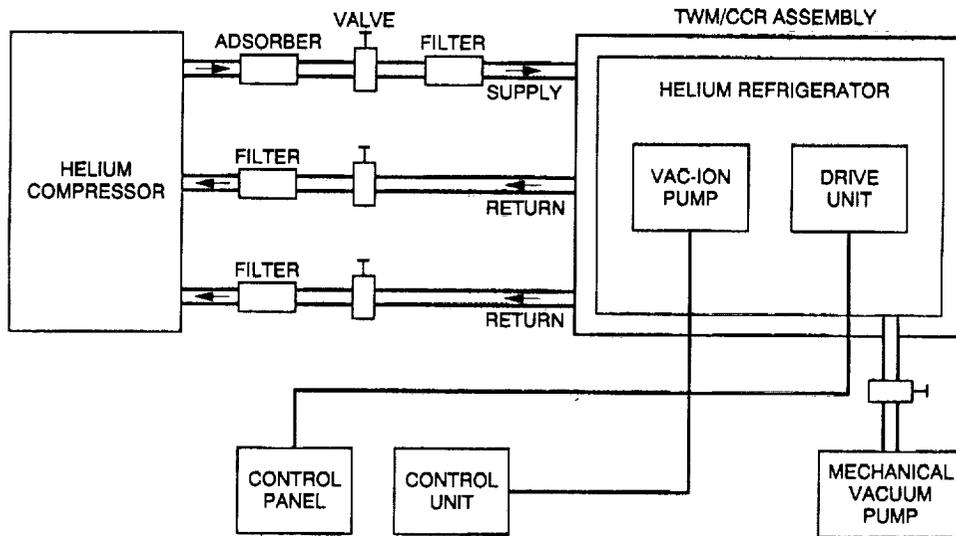


Fig. 2. Simplified TWM/CCR Group Equipment Connections diagram.

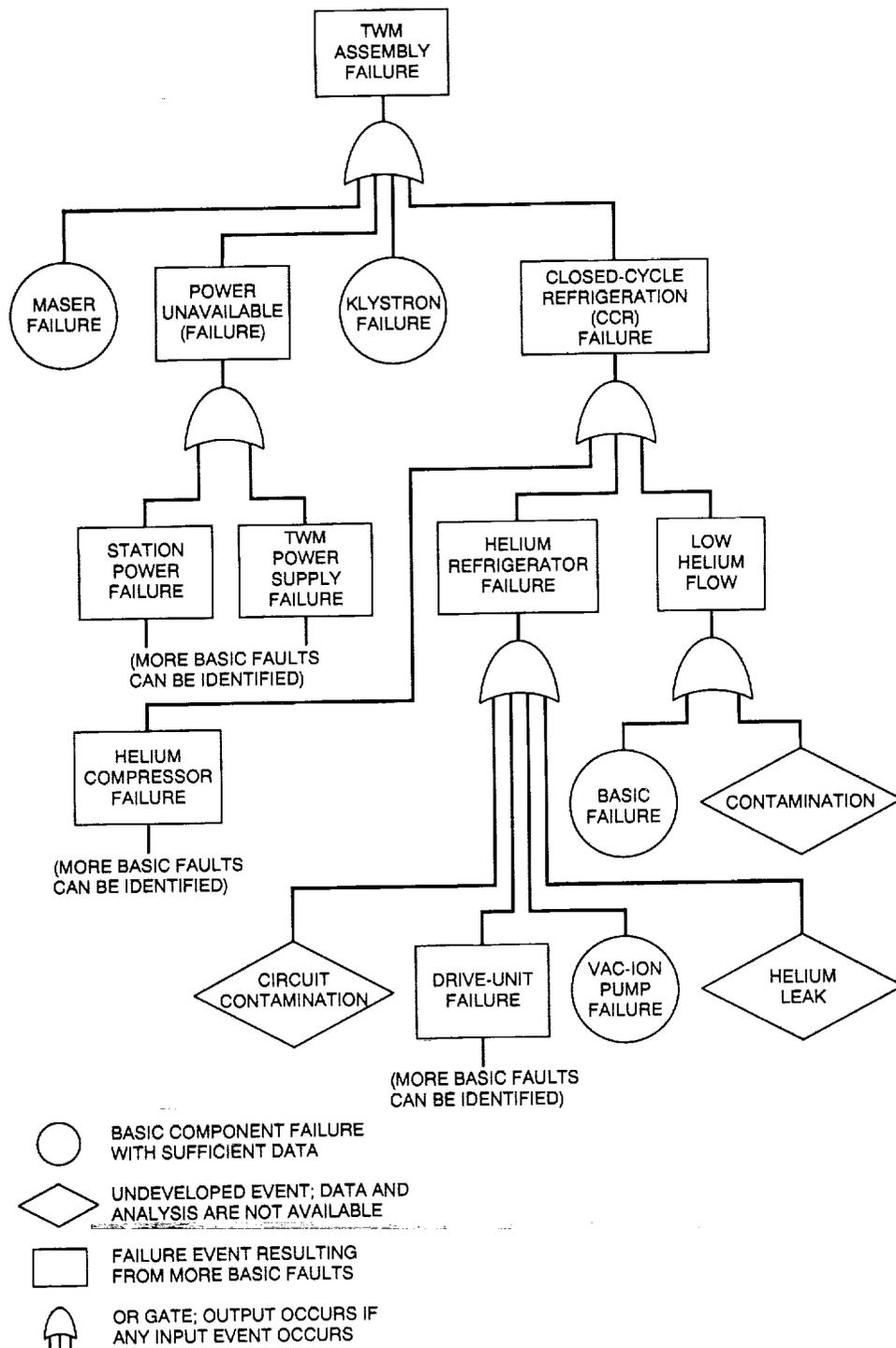
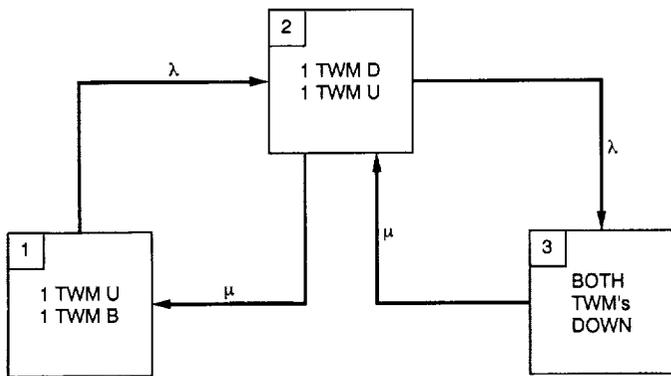


Fig. 3. Fault tree for a TWM assembly (top failure events are shown; basic component failures are not specified).



U UP, OPERATING
D DOWN
B UP, BACKUP

Fig. 4. State-space performance model for a two-assembly standby TWM configuration.

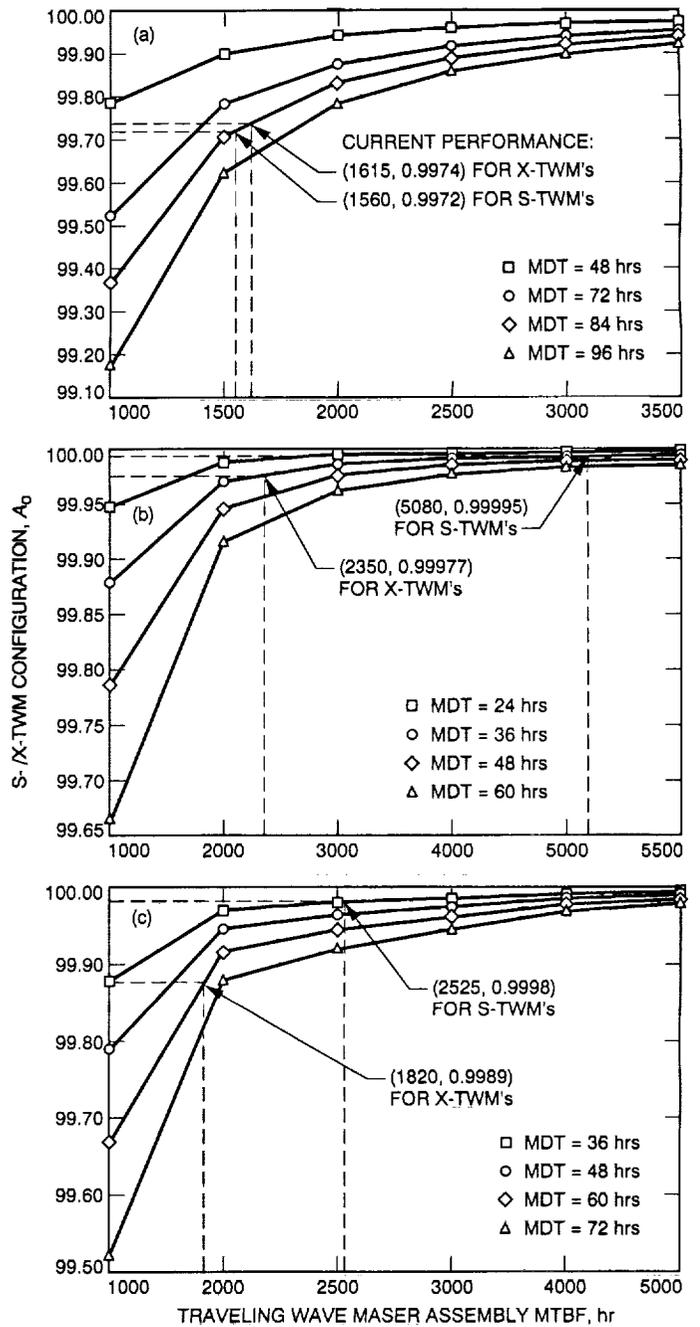


Fig. 5. Operational availability of S- and X-band TWM configuration versus TWM assembly MTBF at various levels of TWM Assembly MDT: (a) Goldstone, DSS 14; (b) Canberra, DSS 43; and (c) Madrid, DSS 63.

523-83

128493-19436

A Model for the Cost of Doing a Cost Estimate

D. S. Remer¹

Harvey Mudd College of Engineering and Science
Claremont, California

H. R. Buchanan

Radio Frequency and Microwave Subsystems Section

This article provides a model for estimating the cost required to do a cost estimate for DSN projects that range from \$0.1 to \$100 million. The cost of the cost estimate in thousands of dollars, C_E , is found to be approximately given by $C_E = KC_p^{0.35}$, where C_p is the cost of the project being estimated in millions of dollars and K is a constant depending on the accuracy of the estimate. For an order-of-magnitude estimate, $K = 24$; for a budget estimate, $K = 60$; and for a definitive estimate, $K = 115$. That is, for a specific project, the cost of doing a budget estimate is about 2.5 times as much as that for an order-of-magnitude estimate, and a definitive estimate costs about twice as much as a budget estimate. Use of this model should help provide the level of resources required for doing cost estimates and, as a result, provide insights towards more accurate estimates with less potential for cost overruns.

I. Introduction

Large cost overruns for major projects are a frequent occurrence. For example, the following projects are reported to have had final costs that exceeded the original cost estimates by over 45 percent.² These were Landsat-D (48 percent), Infrared Astronomical Satellite (60 percent), Earth Radiation Budget Experiment (61 percent), Gamma

Ray Observer (98 percent), Space Telescope (98 percent), Galileo (100 percent), Tracking and Data Relay Satellite System (130 percent), and Pegasus (700 percent). Considering the current emphasis on fiscal responsibility within NASA and other government agencies, cost overruns are a major problem. Overruns may lead to cancellation of the project. In some cases, a potential overrun results in modifying the project to a design-to-cost task.

There are many reasons for cost overruns, but one of the key factors is the lack of resources (time, money, and staffing) spent to do proper up-front cost estimates. An-

¹ Consultant to the TDA Planning Section.

² H. W. Partma and W. E. Ruhland, *Predicting Financial Risk for the Development of Space Flight Projects* (internal document), Jet Propulsion Laboratory, Pasadena, California, September 1988.

other major reason is that the implementers did not do the estimating. The purpose of this article is to address the issue of the cost to do a cost estimate. The authors will report on how others handle this issue, offer suggestions on how the DSN should estimate the amount to spend on a cost estimate, and discuss its impact on reducing the probability of a cost overrun.

The authors will report on their literature search and actual data from JPL Procurement on what others charge the Laboratory for a cost estimate. The goal is to determine guidelines and a methodology for estimating how much to spend on a cost estimate to achieve a desired accuracy. Underallocation of resources for producing a cost estimate is not uncommon. All the necessary cost elements are usually not included because of time constraints. This leads to cost overruns and/or descopeing of the functional requirements of projects.

II. The Cost of Estimating Accuracy

The cost of doing a cost estimate depends on how well the project is defined, who is doing the estimate, the amount of information available, and the level of accuracy required. An order-of-magnitude estimate will cost much less than a definitive estimate. The accuracy of a cost estimate increases within certain limits as the amount of resources spent on the cost estimate increases. The authors defined a metric for the cost to do a cost estimate as the percent of the cost of the estimate as compared with the total cost of the project.

cost of a cost estimate (percent) =

$$\frac{\text{cost of the estimate } (C_E)}{\text{total cost of the project } (C_P)} \times 100 \text{ percent}$$

Figure 1 shows the relative cost of a cost estimate [1] as a function of the accuracy of an estimate for a project costing approximately \$3 million that the authors updated to 1990 dollars using the *NASA Inflation Index* [3]. For example, an estimate that is accurate enough to be within ± 30 percent would cost 0.2 percent, or \$6 thousand, whereas an estimate accurate to within ± 10 percent would cost 1.5 percent, or \$45 thousand, of the total project cost. The more one invests up front in defining the requirements and the deliverables, the more accurate the final estimate will be.

For projects much larger than \$3 million, the cost of the estimate as a percent of the total project cost would be less

than that shown in Fig. 1, whereas for smaller projects costing much less than \$3 million, the percent spent on the cost estimate would be higher. Figure 1 represents a model typical of the process industry; however, the concept applies to the DSN. The authors will now report on a recent set of data that is applicable to the DSN.

This second set of data [2] shows the cost to prepare cost estimates for three accuracy ranges varying from order of magnitude, -30 to $+50$ percent; budget, -15 to $+30$ percent; and definitive, -5 to $+15$ percent, for projects ranging from approximately \$0.1 to \$80 million. Notice that the high limits of the ranges are greater than the low limits because there is usually a lack of consideration of all the necessary cost elements. As a result, there is usually more chance of a cost overrun than an underrun. By making several smoothing assumptions and updating the data to 1990 using the *NASA Inflation Index* [3], the authors plotted the resulting data set as shown on a log-log plot in Fig. 2. A model developed based on these parameters is described below.

III. Model for the Cost of Estimating Accuracy

On the log-log plot of Fig. 2, a set of straight lines conformed closely to the data points. On a log-log plot, a straight line represents a convenient power function equation of the form $C_E = K C_P^R$. That is, by taking the log of both sides of the equation, one gets

$$\log C_E = R \log C_P + \log K \quad (1)$$

This represents a straight line where R is the slope of the line in Fig. 2 and $\log K$ is the Y intercept. The lines shown in Fig. 2 therefore reflect a convenient power function equation that can be used as a model.

$$C_E = K C_P^R \quad (2)$$

where

C_E = cost of the cost estimate in thousands of dollars

C_P = cost of the project being estimated in millions of dollars

K = a constant depending on the accuracy of the estimate

R = slope of lines

Figure 2 shows the slope R and the constant K for each class of cost estimate. For each class of estimate, $R = 0.35$ for project costs in the range of \$0.1 to \$100 million. The constant K is 24 for an order-of-magnitude estimate, 60 for a budget estimate, and 115 for a definitive estimate. Or to look at it another way, a budget cost estimate costs about two-and-one-half times as much as an order-of-magnitude estimate, and a definitive estimate costs about twice as much as a budget estimate.

IV. Discussion of Levels of Cost Estimates

The levels of cost estimates discussed in this article correlate with the condensed classification of cost estimates proposed by the American Association of Cost Engineers [2]. These are as follows:

Class	Accuracy, percent
Order of magnitude	-30 to +50
Budget	-15 to +30
Definitive	-5 to +15

An order-of-magnitude level of cost estimate is usually based on very preliminary statements of requirements. This is done in the requirements definition stage when there is a preliminary listing of deliverables. This class of estimate roughly coincides with that needed for a Level A design review³ when a maximum uncertainty of 30 percent is desired.

The budgetary level of a cost estimate is based on system functional requirements with at least preliminary deliverables, receivables, and schedules presented by a subsystem. This class of cost estimate is appropriate for Level B and/or Level C design reviews when a maximum uncertainty of 20 percent is desired.

The definitive level of a cost estimate is based on a subsystem functional design, and the deliverables, receivables, and schedules are carefully defined and final. This class of cost estimate is appropriate for a Level D design review with a maximum uncertainty of 10 percent.

³ R. P. Mathison and P. T. Westmoreland, "Cost Review Format," JPL Interoffice Memorandum 3300-88-08 (internal document), Jet Propulsion Laboratory, Pasadena, California, January 6, 1988.

A more detailed description of the DSN classes of cost estimates as they relate to design reviews is presented in the Mathison-Westmoreland JPL Interoffice Memorandum.⁴

V. Example Using the Model

Assume that one has to estimate the cost required to do a cost estimate for a project that is expected, based on other similar projects, to cost approximately \$20 million. Use Eq. (2) or Fig. 2 where

$$C_E = KC_P^R$$

$C_P = 20$, $R = 0.35$, and $K = 24$, 60, and 115 for an order-of-magnitude, a budget, and a definitive estimate, respectively. Using $C_E = 24 \times 20^{0.35}$, one gets \$68,000 for an order-of-magnitude estimate. For a budget estimate, one gets \$171,000 and a definitive estimate costs \$328,000. Armed with these data, a decision can be made to proceed with the cost estimate after allocating the necessary funds.

This method may reduce underallocation of resources for producing cost estimates, and thereby more realistic project cost estimates may be obtained. Of course, if the actual estimate of the project turns out to be more or less than the so-called ballpark guesstimate, the budget for doing the cost estimate can be adjusted accordingly. In the next section, data obtained from JPL Procurement will be presented on the cost of actual cost estimates.

VI. JPL Procurement Data for Cost Estimates

The authors obtained data based on JPL procurements for outside contractors to do cost estimates for DSN projects.^{5,6,7} These data points are summarized in Table 1. The first data point reflects a Motorola estimate

⁴ Ibid.

⁵ R. S. Hughes, personal communication, Radio Frequency and Microwave Subsystems Section, Jet Propulsion Laboratory, Pasadena, California, December 5, 1989.

⁶ R. L. White, personal communication, Ground Antennas and Facilities Engineering Section, Jet Propulsion Laboratory, Pasadena, California, January 23, 1991.

⁷ L. H. Kushner, personal communication, Ground Antennas and Facilities Engineering Section, Jet Propulsion Laboratory, Pasadena, California, February 5, 1991.

for a significant supplement to an existing Motorola contract. The second and third data points show the costs of two externally generated Preliminary Engineering Reports (PER's) that developed the estimated costs for cost-of-facilities projects

The costs of the cost estimates for these three projects varied from 1.3 to 2.9 percent of the total project cost. The high value of 2.9 percent was for a relatively small project of about \$2 million, whereas the lower values of 1.3 to 1.5 percent were for projects in the \$11 to \$24 million range. These results fall into the band of curves shown in

Fig. 2. This provides an independent check of the model proposed earlier.

VII. Summary

A model for estimating how much should be allocated to do DSN cost estimates for new capabilities has been developed. This model may help the DSN make better cost estimates and thereby reduce the possibility of producing cost estimates that are too low. These low cost estimates could lead to cost overruns or reduction of some functional requirements or both.

References

- [1] W. R. Park, *Cost Engineering Analysis—A Guide to the Economic Evolution of Engineering Projects*, New York: John Wiley & Sons, p. 133, 1973.
- [2] K. K. Humphreys, *Jelen's Cost and Optimization Engineering*, third edition, New York: McGraw-Hill, p. 370, 1991.
- [3] National Aeronautics and Space Administration, Cost and Economic Analysis Branch, *NASA Inflation Index*, Washington, DC, February 28, 1991.

Table 1. JPL procured cost estimates in FY'90 dollars.

Source	C_P , millions of dollars	C_E , thousands of dollars	$(C_E/C_P)100$, percent
Motorola ^a	1.96	56	2.9
Section 332 ^b PER ^c	10.83	146	1.3
Section 332 PER ^d	24.00	360	1.5

^a Modification to Motorola contract (Magellan ground hardware) for adding C-band uplink capability to DSN receiver-exciter subsystems.

^b Ground Antennas and Facilities Engineering Section.

^c For 34-m antenna JPL support effort plus contractor production of PER.

^d For new Telecommunication Research Laboratory (building).

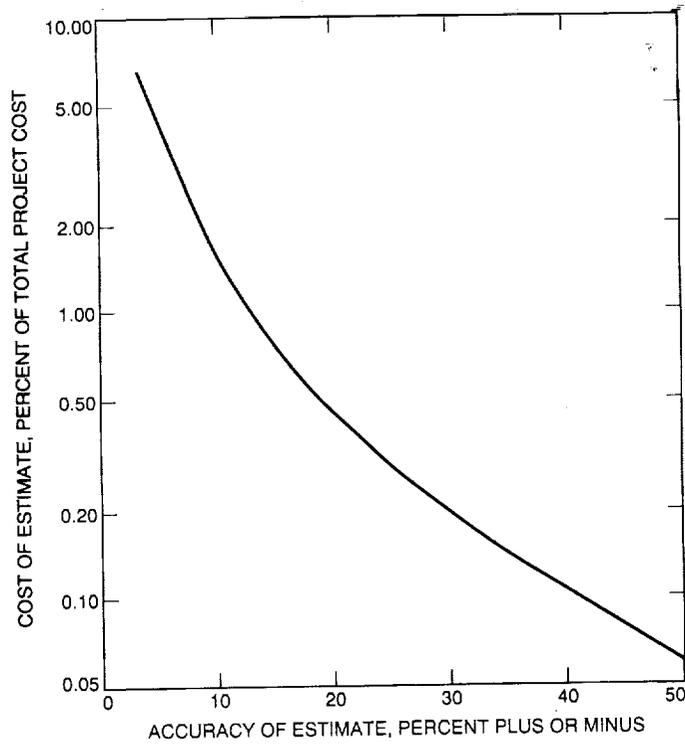


Fig. 1. Cost of a cost estimate for a \$3 million (1990) project [1].

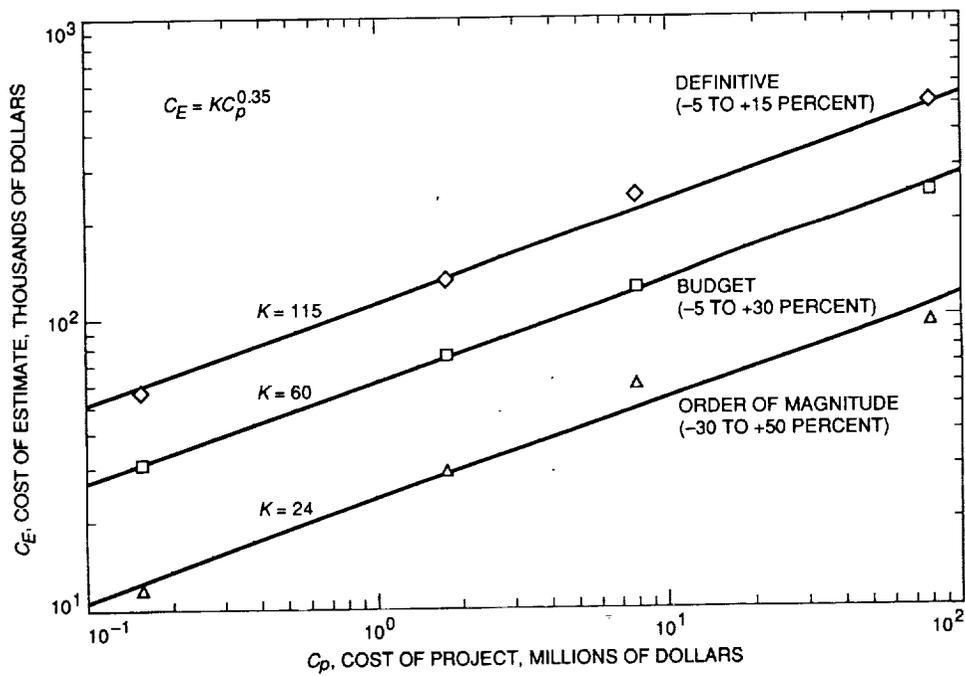


Fig. 2. Cost of a cost estimate for three accuracy ranges.

Referees

The following people have refereed articles for *The Telecommunications and Data Acquisition Progress Report*. By attesting to the technical and archival value of the articles, they have helped to maintain the excellence of this publication during the past year.

H. Ansari	G. J. Dick	J. W. Layland	M. Simon
D. A. Bathker	J. O. Dickey	F. Manshadi	R. Syndor
J. S. Border	R. Dickinson	R. J. McEliece	R. C. Tausworthe
D. Boussalis	S. Dolinar	F. Pollara	S. W. Thurman
S. Butman	W. Gawronski	R. Rasmussen	J. S. Ulvestad
A. Cha	R. A. Jacobson	E. H. Satorius	C. L. Zygielbaum
P. A. Clements	A. S. Konopliv	J. Shell	

ORIGINAL PAGE IS
OF POOR QUALITY